

Testing linearity in semiparametric regression models

TAPIO NUMMI*, JIANXIN PAN AND NICHOLAS MESUE

One of the fundamental assumptions of a basic multiple linear regression model is that the contribution of each of the model terms is strictly linear. In many cases, this may be an excessive simplification of the complicated relationships. Moreover, it may be difficult or impossible to test the hypothesized model against all possible kinds of relevant alternative models. Therefore, tests that perform well under more general circumstances are also required. This paper considers the semiparametric model, where the contribution of one of the model terms may not be strictly linear, and also proposes an exact F-test for the situation. The method also allows dependent error terms. The performance of the proposed test is illustrated by simulation experiments and in real air pollution and health data.

AMS 2000 SUBJECT CLASSIFICATIONS: Primary 62G08, 60K35; secondary 60K35.

KEYWORDS AND PHRASES: Cubic smoothing splines, Fine particles, F-test, Partial linear model.

1. INTRODUCTION

In many scientific fields, it is often necessary to establish a model that accurately describes real-life relationships. Essentially, the modeling challenge can be driven by the underlying task or data, or both. A traditional statistical approach to modeling is based on linear or non-linear regression model and least squares (or maximum likelihood) methods. Generalized linear models can be considered to be a further extension of these methods. The basic feature in these models is that they are based on a strict assumption of a parametric model. In some cases, these approaches may provide a nice presentation of the data at hand. In many cases, however, the true model form is not known and model choice is restricted to rather arbitrary, parametric competitive models. Over the last two decades, nonparametric methods such as Kernel estimation, local polynomials, regression splines, smoothing splines, and generalized additive models have provided a serious alternative to parametric statistical methods (see e.g. [8, 14, 18] and [19]).

The present paper considers an interesting mixture of parametric and nonparametric models. The model considered here is basically a multiple linear regression model with

*Corresponding author.

a possible smooth term. The study was primarily interested in testing whether the smooth term could be exchanged to the linear term so that the full multiple regression model could be used instead of a semiparametric model. Statistical inference of semiparametric regression models have been considered by several authors. However, most of the references concentrate on the methods for independent data (see e.g [1, 2, 5] and [7] for linear models and [10] and [11] for generalized linear models). For dependent data, we may refer to [20] who have considered statistical inference under the semiparametric additive model framework. In summary, most of the references rely on the independence assumption and the exact distribution of the test statistics is either not known or is fairly complicated. Our test relies on the exact F-distribution and also applies under correlated observations. The test is computationally very simple and it is not influenced by the correlation of the explanatory variables and the smooth term.

This article is structured as follows. Section 2 presents the basic cubic smoothing spline model and its estimation. Section 3 extends these results to a semiparametric model and proposes a new F-test based on a semiparametric model. In Section 4 the performance of the proposed test is illustrated by simulation experiments and in real air pollution and health data from the city of Tampere, Finland. Concluding remarks are provided in Section 5.

2. CUBIC SMOOTHING SPLINES

Suppose that $\mathbf{y} = (y_1, y_2, \dots, y_n)'$ is the $n \times 1$ vector of observations and that design points (knots) x_1, x_2, \dots, x_n are given at a certain interval $[a, b]$ satisfying $a < x_1 < x_2 < \dots < x_n < b$. A cubic smoothing spline can now be written as follows

$$(1) \quad \mathbf{y} = \mathbf{g} + \epsilon,$$

where $\mathbf{g} = (g(x_1), g(x_2), \dots, g(x_n))'$, $g(\cdot)$ is a smooth twice differentiable curve and $\epsilon = (\epsilon_1, \dots, \epsilon_n)' \sim N(\mathbf{0}, \sigma^2 \mathbf{R})$, where \mathbf{R} is a covariance matrix with parameters θ .

Let $h_j = x_{j+1} - x_j$, $j = 1, 2, \dots, (n - 1)$. Define the non-zero elements of banded $n \times (n - 2)$ and $(n - 2) \times (n - 2)$ matrices ∇ and Δ , respectively, as

$$\nabla_{k,k} = \frac{1}{h_k}, \quad \nabla_{k+1,k} = -\left(\frac{1}{h_k} + \frac{1}{h_{k+1}}\right), \quad \nabla_{k+2,k} = \frac{1}{h_{k+1}}$$

1 and

$$2 \quad \Delta_{k,k+1} = \Delta_{k+1,k} = \frac{h_{k+1}}{6}, \quad \Delta_{k,k} = \frac{h_k + h_{k+1}}{3},$$

3 where $k = 1, 2, \dots, (n-2)$. If the roughness matrix is de-
 4 noted as $\mathbf{K} = \nabla \mathbf{\Delta}^{-1} \nabla'$, the solution (at the design points)
 5 to the penalized least squares criterion (PLS)

$$6 \quad (2) \quad Q_1 = (\mathbf{y} - \mathbf{g})' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{g}) + \alpha \mathbf{g}' \mathbf{K} \mathbf{g}$$

7 for fixed positive α is a cubic smoothing spline

$$8 \quad (3) \quad \tilde{\mathbf{g}} = (\mathbf{H} + \alpha \mathbf{K})^{-1} \mathbf{H} \mathbf{y} = \mathbf{S}_\alpha \mathbf{y},$$

9 where it is denoted $\mathbf{R}^{-1} = \mathbf{H}$ and the so-called *smoother*
 10 *matrix* is $\mathbf{S}_\alpha = (\mathbf{H} + \alpha \mathbf{K})^{-1} \mathbf{H}$. If the covariance matrix \mathbf{R}
 11 satisfies the equation $\mathbf{K} \mathbf{R} = \mathbf{K}$ then

$$12 \quad (4) \quad \mathbf{K} = \mathbf{K} \mathbf{H}$$

13 and the smoother matrix reduces to the simplified form
 14 $\mathbf{S}_\alpha = (\mathbf{I} + \alpha \mathbf{K})^{-1}$, which does not explicitly involve covari-
 15 ance matrix \mathbf{R} . The simplified estimator is then

$$16 \quad (5) \quad \hat{\mathbf{g}} = (\mathbf{I} + \alpha \mathbf{K})^{-1} \mathbf{y}.$$

17 Note that, in this case, spline estimates are simple linear
 18 functions of the observations y_1, y_2, \dots, y_n . A set of covari-
 19 ance matrices satisfying the condition (4) can be generated
 20 using the formulas discussed in [12] and [13]. As a special
 21 case these structures include, for example, $\mathbf{R} = \mathbf{X} \mathbf{D} \mathbf{X}' + \mathbf{I}$,
 22 where $\mathbf{X} = [\mathbf{1}, \mathbf{x}]$, $\mathbf{1} = (1, \dots, 1)'$ and $\mathbf{x} = (x_1, \dots, x_n)'$. Note
 23 that for simplicity, hereinafter we take $\mathbf{R} = \mathbf{X} \mathbf{D} \mathbf{X}' + \mathbf{I}$, but
 24 our results also applies under more general class of covari-
 25 ance structures discussed in [13].

26 Generally, the smoother matrix $\mathbf{S}_\alpha = (\mathbf{I} + \alpha \mathbf{K})^{-1}$ is not
 27 a projection matrix; therefore, for example, the theory of
 28 linear models is not directly applicable for statistical infer-
 29 ence of smoothing splines. However, it is possible to plausi-
 30 bly approximate \mathbf{S}_α using a projection matrix [13]. Let the
 31 eigenvalue decomposition of roughness matrix \mathbf{K} be

$$32 \quad (6) \quad \mathbf{K} = \mathbf{T} \mathbf{\Lambda} \mathbf{T}',$$

33 where \mathbf{T} is the matrix of n orthonormal eigenvectors and $\mathbf{\Lambda}$
 34 is the diagonal matrix of eigenvalues $\lambda_1, \dots, \lambda_n$ of \mathbf{K} . It can
 35 easily be seen that

$$36 \quad (7) \quad \mathbf{S}_\alpha = \mathbf{T} (\mathbf{I} + \alpha \mathbf{\Lambda})^{-1} \mathbf{T}'.$$

37 Note that the set of eigenvectors of \mathbf{S}_α and \mathbf{K} is the same
 38 and the eigenvalues are connected such that the eigenvalues
 39 of \mathbf{S}_α are $\gamma = 1/(1 + \alpha \lambda)$. Hereinafter we assume that eigen-
 40 vectors are ordered according to eigenvalues of \mathbf{S}_α . The first
 41 two eigenvalues of \mathbf{S}_α are always one and the first two eigen-
 42 vectors span the subspace corresponding to the straight line
 43 model [9]. In general the sequence of eigenvectors $\mathbf{t}_1, \dots, \mathbf{t}_n$

57 appears to increase in complexity like a sequence of orthog-
 58 onal polynomials defined on \mathbf{x} . The spline approximation
 59 used here is based on the $c > 2$ first eigenvectors of \mathbf{T} . This
 60 approximation minimizes the least squares criterion

$$61 \quad (\mathbf{y} - \mathbf{T}_* \xi)' (\mathbf{y} - \mathbf{T}_* \xi),$$

62 where \mathbf{T}_* is the matrix of c first eigenvectors of \mathbf{T} and the
 63 fit is simply

$$64 \quad \tilde{\mathbf{y}} = \mathbf{P} \mathbf{y},$$

65 where $\mathbf{P} = \mathbf{T}_* \mathbf{T}_*'$. It was shown in [13] that the approxima-
 66 tion is good for relatively smooth data. It is computationally
 67 simple (calculated directly from \mathbf{K}) and \mathbf{P} is now a projec-
 68 tion matrix, which makes the application of the theory of
 69 linear models possible in this framework.

70 3. SEMIPARAMETRIC MODEL

71 3.1 Estimation

72 The semiparametric model to be considered here is an
 73 extension of the basic spline model of previous section. The
 74 model takes the following form

$$75 \quad (8) \quad \mathbf{y} = \mathbf{U} \mathbf{b} + \mathbf{g} + \epsilon$$

76 where the extension is the linear part $\mathbf{U} \mathbf{b}$, where \mathbf{U} is a full
 77 rank $n \times k$ model matrix of k explanatory variables (constant
 78 term not included) and \mathbf{b} is a k -vector of unknown param-
 79 eters. Estimation of semiparametric models have been con-
 80 sidered in [4, 15] and [17] and applied in [16], for example.
 81 Our approach is to minimize the penalized sum of squares
 82 (see [4])

$$83 \quad (9) \quad Q_2 = [\mathbf{y} - (\mathbf{U} \mathbf{b} + \mathbf{g})]' \mathbf{H} [\mathbf{y} - (\mathbf{U} \mathbf{b} + \mathbf{g})] + \alpha \mathbf{g}' \mathbf{K} \mathbf{g}.$$

84 The minimization with respect to \mathbf{b} and \mathbf{g} yields to estimates

$$85 \quad (10) \quad \tilde{\mathbf{b}} = [\mathbf{U}' \mathbf{H} (\mathbf{I} - \mathbf{S}_\alpha) \mathbf{U}]^{-1} \mathbf{U}' \mathbf{H} (\mathbf{I} - \mathbf{S}_\alpha) \mathbf{y}$$

86 and

$$87 \quad (11) \quad \tilde{\mathbf{g}} = \mathbf{S}_\alpha (\mathbf{y} - \mathbf{U} \tilde{\mathbf{b}})$$

88 where $\mathbf{S}_\alpha = (\mathbf{H} + \alpha \mathbf{K})^{-1} \mathbf{H}$. If the covariance structure is
 89 assumed to be, for example $\mathbf{R} = \mathbf{X} \mathbf{D} \mathbf{X}' + \mathbf{I}$, the condition
 90 $\mathbf{K} = \mathbf{K} \mathbf{H}$ holds and it can be shown that

$$91 \quad \mathbf{H} (\mathbf{I} - \mathbf{S}_\alpha) = \mathbf{I} - \mathbf{S}_\alpha$$

92 (see Appendix A for details). Therefore, the simplified es-
 93 timators are

$$94 \quad (12) \quad \hat{\mathbf{b}} = [\mathbf{U}' (\mathbf{I} - \mathbf{S}_\alpha) \mathbf{U}]^{-1} \mathbf{U}' (\mathbf{I} - \mathbf{S}_\alpha) \mathbf{y}$$

95 and $\hat{\mathbf{g}} = \mathbf{S}_\alpha (\mathbf{y} - \mathbf{U} \hat{\mathbf{b}})$, where $\mathbf{S}_\alpha = (\mathbf{I} + \alpha \mathbf{K})^{-1}$. The whole
 96 semiparametric curve is then fitted by

$$97 \quad (13) \quad \hat{\mu} = \mathbf{M} \mathbf{y},$$

1 where $\mathbf{M} = \mathbf{S}_\alpha + \tilde{\mathbf{U}}[\tilde{\mathbf{U}}'\mathbf{U}]^{-1}\tilde{\mathbf{U}}'$ and $\tilde{\mathbf{U}} = (\mathbf{I} - \mathbf{S}_\alpha)\mathbf{U}$, respec-
 2 tively. Fixing the dimension of the spline approximation (see
 3 the next section for details) gives

$$4 \quad (14) \quad \hat{\boldsymbol{\mu}} = \mathbf{M}_*\mathbf{y},$$

5 where $\mathbf{M}_* = \mathbf{P} + \mathbf{U}_*(\mathbf{U}'_*\mathbf{U}_*)^{-1}\mathbf{U}'_*$ and $\mathbf{U}_* = (\mathbf{I} - \mathbf{P})\mathbf{U}$.
 6 Note that \mathbf{M}_* is now a projection matrix and thus lays the
 7 grounds for the application of the standard theory of linear
 8 models.

9 Using the approximation (14) makes it easy to investigate
 10 the statistical properties of the estimates of the linear part of
 11 the model. As discussed in the previous section, parameter
 12 estimates for \mathbf{b} are then obtained by

$$13 \quad (15) \quad \hat{\mathbf{b}} = (\mathbf{U}'_*\mathbf{U}_*)^{-1}\mathbf{U}'_*\mathbf{y}$$

14 and, if $\mathbf{R} = \mathbf{XDX}' + \mathbf{I}$ and $\mathbf{U}'_*\mathbf{X} = \mathbf{0}$, it can easily be seen
 15 that

$$16 \quad (16) \quad \text{Var}(\hat{\mathbf{b}}) = \sigma^2(\mathbf{U}'_*\mathbf{U}_*)^{-1}.$$

17 3.2 Testing the semiparametric model

18 This study aims to test whether our model is the full
 19 linear model. The null hypothesis (*linear model*) is as follows

$$20 \quad \mathbf{H}_0 : \mu = \mathbf{U}_{k+2}\mathbf{b}_{k+2},$$

21 where $\mathbf{U}_{k+2} = [\mathbf{U}, \mathbf{X}]$, where $\mathbf{X} = [\mathbf{1}, \mathbf{x}]$ and the alternative
 22 hypothesis (*semiparametric model*) is

$$23 \quad \mathbf{H}_a : \mu = \mathbf{U}\mathbf{b} + \mathbf{g},$$

24 where \mathbf{g} is a smooth term. Testing is now based on the fit
 25 obtained by the approximation (14) with $\hat{\boldsymbol{\mu}} = \mathbf{M}_*\mathbf{y}$. If we
 26 take $\mathbf{R} = \mathbf{XDX}' + \mathbf{I}$, for example, $\mathbf{M}_*(\mathbf{XDX}') = \mathbf{XDX}'$
 27 and it can easily be seen that

$$28 \quad (\mathbf{I} - \mathbf{M}_*)(\mathbf{XDX}' + \mathbf{I})(\mathbf{I} - \mathbf{M}_*) = \mathbf{I} - \mathbf{M}_*$$

29 and therefore the distribution of (see also [3])

$$30 \quad \sigma^{-2}\mathbf{y}'(\mathbf{I} - \mathbf{M}_*)\mathbf{y} = \sigma^{-2}S_{min} \sim \chi_{n-c-k}^2,$$

31 where we denote $\mathbf{y}'(\mathbf{I} - \mathbf{M}_*)\mathbf{y} = S_{min}$. If now in \mathbf{M}_* we
 32 define $\mathbf{P} = \mathbf{T}_*\mathbf{T}'_*$ where $\mathbf{T}_* = [\mathbf{t}_1, \mathbf{t}_2]$ (vectors spanning the
 33 linear subspace) and denote the projection matrix as \mathbf{M}_{OLS} .
 34 Then the distribution of

$$35 \quad \sigma^{-2}\mathbf{y}'(\mathbf{I} - \mathbf{M}_{OLS})\mathbf{y} = \sigma^{-2}S_{OLS} \sim \chi_{n-k-2}^2,$$

36 where we denote $\mathbf{y}'(\mathbf{I} - \mathbf{M}_{OLS})\mathbf{y} = S_{OLS}$, and further

$$37 \quad \sigma^{-2}(S_{OLS} - S_{min}) \sim \chi_{c-2}^2.$$

38 Since $\mathbf{M}_*\mathbf{M}_{OLS} = \mathbf{M}_{OLS}\mathbf{M}_* = \mathbf{M}_{OLS}$ it can also be easily
 39 observed that $(\mathbf{I} - \mathbf{M}_*)(\mathbf{M}_* - \mathbf{M}_{OLS}) = \mathbf{0}$ and therefore

40 the sum of squares S_{min} and $S_{OLS} - S_{min}$ are independent.
 41 Testing can now be based on

$$42 \quad (17) \quad F = \frac{(S_{OLS} - S_{min})/(c-2)}{S_{min}/(n-c-k)} \sim F(c-2, n-c-k).$$

43 Then observing the larger value for F than the quantile
 44 $F_{1-\alpha}(c-2, n-c-k)$ yields rejection of the null hypothesis
 45 of the full linear model. Note that since \mathbf{P} and \mathbf{U}_* in \mathbf{M}_*
 46 are orthogonal F is not influenced by the correlation of the
 47 explanatory variables \mathbf{U} and the smooth term \mathbf{g} . It should
 48 also be emphasized that the exact results holds for fixed c
 49 and for certain forms of correlation, e.g. $\mathbf{R} = \mathbf{XDX}' + \mathbf{I}$. In
 50 the next section, empirical power of the proposed test was
 51 studied when c is estimated from the data.

52 4. NUMERICAL EXAMPLES

53 4.1 A simulation study

54 To investigate the empirical power of our methods we
 55 conducted a simulation study. We first introduce the mixed
 56 model formulation

$$57 \quad \hat{\mathbf{g}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\boldsymbol{\xi}}$$

58 of the spline solution (3). The fitted spline is then obtained
 59 as the BLUE (best linear unbiased) and BLUP (best linear
 60 unbiased predictor) solutions of the mixed model (see [13])

$$61 \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\xi} + \boldsymbol{\epsilon}$$

62 where $\mathbf{Z} = \nabla(\nabla'\nabla)^{-1}\boldsymbol{\Delta}^{-1/2}$, $\boldsymbol{\xi}$ and $\boldsymbol{\epsilon}$ are independently dis-
 63 tributed as $\boldsymbol{\xi} \sim N(\mathbf{0}, \sigma_\xi^2\mathbf{I})$ and $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$. The smoothing
 64 parameter can be viewed as the ratio $\alpha = \sigma^2\sigma_u^2$. If we define
 65 the effective degrees of freedom as $edf = \text{tr}(\mathbf{S}_\alpha)$ we can solve
 66 for fixed values of edf and σ^2 the corresponding value of σ_ξ^2 .
 67 In our simulations random vectors were generated from the
 68 model

$$69 \quad (18) \quad \mathbf{y} = \mathbf{U}\mathbf{b} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\xi} + \boldsymbol{\epsilon},$$

70 where $\mathbf{U}\mathbf{b}$ is the linear part and the smooth term is approx-
 71 imated by $\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\xi}$.

72 For our study 1000 random vectors were generated
 73 from (18) for each value of $edf = 2, 4, 6, 8, 10, 15, 20, 25, 30$
 74 with the corresponding σ_ξ^2 . We took $\mathbf{X} = (\mathbf{1}, \mathbf{x})$, $\mathbf{x} =$
 75 $(1, 2, \dots, 100)'$, $\boldsymbol{\beta} = (500, 1)'$ and $\sigma^2 = 15$. The linear
 76 part $\mathbf{U}\mathbf{b}$ of the model was defined using the column vec-
 77 tor $\mathbf{u} = \mathbf{x} + \mathbf{e}$ with $\mathbf{b} = 1$, where \mathbf{e} is the vector of 100
 78 realizations of independent standard normal variables. This
 79 will create highly correlated smooth and linear parts.

80 Three different methods for the selection of c was used.
 81 The first method is to set c equal to edf (and to correspond-
 82 ing σ_ξ^2) used in simulations. The second method is to mini-
 83 mize the information criterion BIC from $\mathbf{y} = \mathbf{U}\mathbf{b} + \mathbf{T}_*\boldsymbol{\xi} + \boldsymbol{\epsilon}$,
 84 where \mathbf{T}_* contains the c first eigenvectors of \mathbf{S}_α . In the third
 85 method we simply fixed $c = 3$.

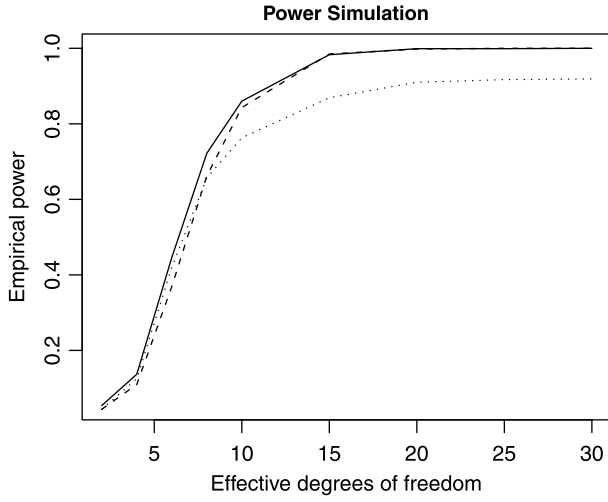


Figure 1. Empirical power as the function of the effective degrees of freedom (dashed line for method 1, solid line for method 2 and dotted line for method 3).

From the Figure 1, it is easily observed that methods 1 and 2 performed practically equally. It seems that basically only minor effects to the empirical power is observed when c is estimated from the data (method 2). These results are also in line with the earlier results by Nummi et al. [13] of the more simple spline model. Also the method 3 performed surprisingly well for low values of $edf < 10$, but for higher values some power is clearly lost.

4.2 Air quality and health in tampere

The example used health and air quality data measured in the medium-sized Finnish city of Tampere daily from January 2006 to August 2008. This data measured air quality variables, weather conditions, traffic density, and health outcomes (such as heart and respiratory diagnosis counts in health centers). This example has concentrated on modeling the fine particle number counts ($< 0, 01 \mu\text{m}^3$). Traffic is known to be one of the main sources of air pollution in urban areas. Heating (connected with temperature) is arguably another important source of air pollution, while humidity has also been found to influence fine particle concentrations. Therefore, these variables were included in the model as explanatory variables. Figure 1 plots a logarithm of fine particle counts with temperature. On the basis of Figure 2, it can be argued that the effects of temperature is not linear within the range of observed values, whereas the effects of traffic (mean number of vehicles/h) and humidity (relative %) are assumed to be linear in the present model. The model for the log fine particle counts is as follows

$$\mathbf{y} = \mathbf{U}\mathbf{b} + \mathbf{g} + \epsilon,$$

where the columns of \mathbf{U} consist of measurements for traffic and humidity and \mathbf{g} represents the effects of temperature.

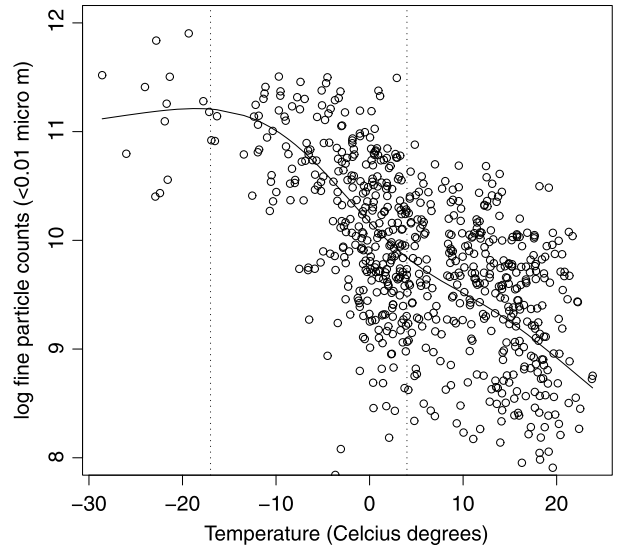


Figure 2. Log fine particle counts as a function of temperature. Note: Some of the original data was excluded due to missing values.

Thus far, we have assumed separate measuring points x_1, \dots, x_n for the smooth term \mathbf{g} for the model. Unfortunately, the data has multiple measuring times; therefore, instead of the smoother $\mathbf{S}_\alpha = (\mathbf{I} + \alpha\mathbf{K})^{-1}$ we used

$$\mathbf{S}_\alpha = \mathbf{N}(\mathbf{N}\mathbf{N}' + \alpha\mathbf{K})^{-1}\mathbf{N}',$$

where \mathbf{N} is an incidence matrix of measuring times. The smoother matrix \mathbf{S}_α was approximated by $\mathbf{T}_* \mathbf{T}_*'$, where the dimension of the approximation was chosen to minimize the BIC criterion. The values of BIC for $c = 1, 2, \dots, 6$ were: 2518.464, 1205.402, 1211.749, 1207.847, 1205.265 and 1211.364. The minimum is obtained at $c = 5$; therefore, this was chosen as the dimension of the approximation. For this model $S_{min} = 214.3947$ and $S_{OLS} = 220.6037$. This gives

$$F = \frac{6.20904/3}{214.3947/685} = 6.612715,$$

which is clearly greater than the quantile $F_{0.95}(3, 685) = 2.617906$. Therefore, the null hypothesis of a linear regression model is rejected and the alternative hypothesis about semiparametric regression is accepted.

Figure 2 plots the fitted smooth term of the semiparametric model with linear terms fixed to their mean values. The effects of temperature on the log fine particles counts is clearly not linear (when also other predictors are included). The effects of temperature increase dramatically below approximately 4 Celcius, but no increase in fine particles counts are observed below -18 Celcius. According to the plot of residuals (Figure 3) the model fits the data very well.

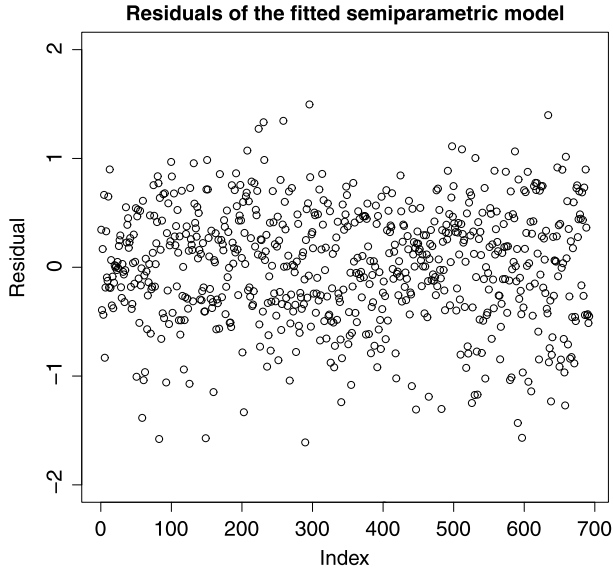


Figure 3. Residuals of the fitted semiparametric model.

5. CONCLUDING REMARKS

This paper has presented a technique that is useful for investigating whether the contribution of a certain covariate in a basic linear model is actually linear or whether it is better to resort to semiparametric modeling. The approach can be useful in situations where the form of the contribution of a covariate is not known (for example, in many medical applications) or if the relation is distorted by outlying observations. The test itself is exact for fixed dimension of the approximation and for certain kind of correlated sequences. However, estimation of the dimension of the approximation seem to have only a minor influence to the empirical power. The test is computationally very simple and it is not influenced by the correlation of the explanatory variables and the smooth term.

APPENDIX A. APPENDIX SECTION

The result is obtained by applying a well known matrix inversion theorem

$$(\mathbf{A} + \mathbf{UBV})^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{UB}(\mathbf{B} + \mathbf{BVA}^{-1}\mathbf{UB})^{-1}\mathbf{BVA}^{-1}$$

to $\mathbf{H} = \mathbf{R}^{-1} = (\mathbf{XDX}' + \mathbf{I})^{-1}$ and $\mathbf{S}_\alpha = (\mathbf{I} + \alpha\mathbf{K})^{-1} = (\mathbf{I} + \alpha\nabla\mathbf{\Delta}^{-1}\nabla')^{-1}$. This yields

$$(\mathbf{I} + \mathbf{XDX}')^{-1} = \mathbf{I} - \mathbf{XD}(\mathbf{D} + \mathbf{DX}'\mathbf{XD})^{-1}\mathbf{DX}'$$

and

$$(\mathbf{I} + \alpha\nabla\mathbf{\Delta}^{-1}\nabla')^{-1} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \alpha\mathbf{\Delta}^{-1})^{-1}\mathbf{Z}',$$

where $\mathbf{Z} = \nabla(\nabla'\nabla)^{-1}$. Since \mathbf{X} and \mathbf{Z} are orthogonal, we can easily obtain

$$\mathbf{H}(\mathbf{I} - \mathbf{S}_\alpha) = \mathbf{I} - \mathbf{S}_\alpha.$$

ACKNOWLEDGEMENTS

The authors are grateful to environmental inspector Milla Hilli-Lukkarinen from the Environmental Protection Unit at the City of Tampere for providing the data set used in the example. The authors are also grateful to the editors and to the associate editor for the comments that led to the improvements of the paper.

Received 2 December 2011

REFERENCES

- [1] AZZALINI, A., and BOWMAN, A. (1993). On the use of non-parametric regression for checking linear relationships. *Journal of the Royal Statistical Society, Series B*, **55**(2) 549–557. [MR1224417](#)
- [2] CRAINICEAU, C., RUPPERT, D., CLAESKENS, G., and WAND, M. P. (2005). Exact likelihood ratio tests for penalized splines. *Biometrika*, **92**(1) 91–103. [MR2158612](#)
- [3] DEMIDENKO, E. (2004). *Mixed Models: Theory and Applications*. Wiley, Hoboken, New Jersey. [MR2077875](#)
- [4] GREEN, P. J., and SILVERMAN, B. W. (1994). *Nonparametric Regression and Generalized Linear Models*. Chapman and Hall, London. [MR1270012](#)
- [5] COX, D., and KOH, E. (1989). A smoothing spline based test of model adequacy in polynomial regression. *Annals of Inst. of Statist. Math.*, **41**(2) 383–400. [MR1006497](#)
- [6] EUBANK, R. L., and HART, J. D. (1992). Testing goodness-of-fit in regression via order selection criteria. *Annals of Inst. of Statist. Math.*, **20**(3) 1412–1425. [MR1186256](#)
- [7] FAN, J., and HUANG, L. S. (2001). Goodness-of fit tests for parametric regression models. *Journal of the American Statistical Association*, **96** 640–652. [MR1946431](#)
- [8] HASTIE, T., and TIBSHIRANI, R. J. (1990). *Generalized Additive Models*. Chapman and Hall, London. [MR1082147](#)
- [9] HASTIE, T., TIBSHIRANI, R., and FRIEDMAN, J. (2001). *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. Springer, New York. [MR1851606](#)
- [10] HÄRDLE, W., MAMMEN, E., and MÜLLER, M. (1998). Testing parametric versus semiparametric modelling in generalized linear models. *Journal of the American Statistical Association*, **93** 1461–1474. [MR1666641](#)
- [11] LIU, A., MEIRING, W., and WANG, Y. (2004). Testing generalized linear models using smoothing spline methods. *Statistica Sinica*, **15** 235–256. [MR2125730](#)
- [12] NUMMI, T., and KOSKELA, L. (2006). Analysis of growth curve data using cubic smoothing splines. *Journal of Applied Statistics*, **35** 681–691. [MR2516865](#)
- [13] NUMMI, T., PAN, J., SIREN, T., and LIU, K. (2011). Testing for cubic smoothing splines under dependent data. *Biometrics*, **67**(3) 871–875. DOI: 10.1111/j.1541-0420.2010.01537.x. [MR2829261](#)
- [14] RUPPERT, D., WAND, M. P., and CARROLL, R. J. (2003). *Semiparametric Regression*. Cambridge University Press, New York, USA. [MR1998720](#)
- [15] SCHIMEK, M. G. (2000). Estimation and inference in partially linear models with smoothing splines. *Journal of Statistical Planning and Inference*, **91** 522–540. [MR1814799](#)
- [16] SCHMALENSEE, R., and STOKER, T. M. (1999). Household gasoline demand in the United States. *Econometrica*, **67** 645–662.
- [17] SPECKMAN, P. (1988). Kernel smoothing in partial linear models. *Journal of the Royal Statistical Society, Series B*, **50** 413–436. [MR0970977](#)
- [18] WOOD, S. N. (2006). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC, Boca Raton, USA. [MR2206355](#)

- 1 [19] WU, H., and ZHANG, J. T. (2006). *Nonparametric Regression*
 2 *Methods for Longitudinal Data Analysis*. Wiley, Hoboken, New
 3 Jersey. [MR2216899](#)
 4 [20] ZHANG, D., and LIN, X. (2003). Hypothesis testing in semipara-
 5 metric additive mixed models. *Biostatistics*, 4(1) 57–74.

6 Tapio Nummi
 7 School of Health Sciences
 8 FI-33014 University of Tampere
 9 Finland
 10 E-mail address: tapio.nummi@uta.fi

Jianxin Pan
 School of Mathematics
 The University of Manchester, M13 9PL
 UK
 E-mail address: jianxin.pan@manchester.ac.uk
 Nicholas Mesue
 School of Health Sciences
 FI-33014 University of Tampere
 Finland
 E-mail address: mesuenicholas@yahoo.com
 url: www.foo.com

57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112