



# HHS Public Access

Author manuscript

*Comput Stat Data Anal.* Author manuscript; available in PMC 2023 January 19.

Published in final edited form as:

*Comput Stat Data Anal.* 2020 October ; 150: . doi:10.1016/j.csda.2020.106987.

## Comparison of nonlinear curves and surfaces

Shi Zhao<sup>a</sup>, Giorgos Bakoyannis<sup>a</sup>, Spencer Lourens<sup>b</sup>, Wanzhu Tu<sup>a</sup>

<sup>a</sup>Department of Biostatistics, Indiana University Fairbanks School of Public Health and Indiana University School of Medicine, Indianapolis, Indiana 46202, U.S.A.

<sup>b</sup>CliftonLarsonAllen LLP

### Abstract

Estimation of nonlinear curves and surfaces has long been the focus of semiparametric and nonparametric regression analysis. What has been less studied is the comparison of nonlinear functions. In lower-dimensional situations, inference typically involves comparisons of curves and surfaces. The existing comparative procedures are subject to various limitations, and few computational tools have been made available for off-the-shelf use. To address these limitations, two modified testing procedures for nonlinear curve and surface comparisons are proposed. The proposed computational tools are implemented in an R package, with a syntax similar to that of the commonly used model fitting packages. An R Shiny application is provided with an interactive interface for analysts who do not use R. The new tests are consistent against fixed alternative hypotheses. Theoretical details are presented in an appendix. Operating characteristics of the proposed tests are assessed against the existing methods. Applications of the methods are illustrated through real data examples.

### Keywords

Comparison of nonlinear functions; Nonparametric and semiparametric regression; Resampling methods

## 1. Introduction

An essential task in nonparametric and semiparametric regression is to estimate nonlinear functions for depiction of relations between independent and dependent variables. In lower-dimensional situations, the functions are often expressed as smooth curves and surfaces [1]. Various smoothing techniques have been developed for the estimation of nonlinear functions. Commonly used methods include local polynomial models [2], wavelets [3], smoothing splines [4, 5], and various types of penalized regression splines [6, 7, 8, 9]. Most of these estimation methods can be easily implemented in common computational platforms, giving analysts much flexibility for curve and surface estimation. What has been less studied

---

wtu1@iu.edu (Wanzhu Tu).

**Publisher's Disclaimer:** This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

is the inference concerning nonlinear functions, and there is a dearth of computational tools for practical use. A question of general interest is whether a specific nonparametric smoother, when applied to different comparison groups, gives the same function.

To address the question above, we review the existing literature on curve and surface comparisons, and present two  $L_2$ -based testing procedures with related theoretical and numerical justifications. Our procedures are based on B-splines, although the formulation of the test statistics could be extended to other smoothers. We put forward an R package, which can be accessed either directly from within R, or through an interactive R-Shiny interface; the latter allows analysts who do not use R to perform the desired comparisons. To illustrate the use of the proposed methods, we present two real data examples.

## 2. Existing methods for curve and surface comparison

Comparison of smooth curves can be formalized as a test of the following hypothesis  $H_0 : g_1(x) = g_2(x) = \dots = g_I(x), \forall x \in \mathbb{R}$  vs  $H_1 : g_i(x) \neq g_j(x)$  for some  $i, j \in \{1, \dots, I\}$ , where  $i$  and  $j$  indicate different comparison groups. In the situation of  $x \in \mathbb{R}^2$ ,  $g_i(\mathbf{x})$  and  $g_j(\mathbf{x})$  are surface functions. The concept can be extended to higher-dimensional functions, although visualization of higher-dimensional functions becomes more difficult. In this paper, we restrict the discussion to lower-dimensional situations, and we write the underlying model as  $Y = g(x) + \epsilon$ . We use capital letters  $Y$  and  $X$  to indicate the random response and independent variables, and their lower-case counterparts  $x$  and  $y$  to indicate the observed values of the corresponding random variables.

Early work on this problem started almost three decades ago. One approach is to frame the problem in a regression setting, where a modified version of the Kolmogorov-Smirnov test could be used to compare the regression curves [10]. Another approach is to transform the nonparametric curves to reduce the comparison to a test of limited dimensional parameters within the transformation matrix [11]. Alternatively, wavelet methods have been used to compare density functions [12]; the methods are especially suitable for comparing higher frequency local features. Many of these methods, however, require the curves to have the same design points, i.e., all functions must be evaluated at the same  $x$  values. To remedy, Kulasekera (1995) fitted kernel-based regression models and proposed tests based on the weighted average of partial sum of squares of the quasi-residuals and error variances [13]. But simulations suggest that these tests tend to be overly sensitive to bandwidth selection. There are also specialized tests for parallelism among the curves [14, 15].

In this section, we briefly review the methods that are most frequently used in analytical practice.

### 2.1. Nonparametric Analysis of Covariance (ANCOVA)

Young and Bowman (1995)[16] described a method for testing the equality of two or more smooth curves, under the model  $Y_{ij} = g_i(x_{ij}) + \epsilon_{ij}$ , where  $\epsilon_{ij} \sim N(0, \sigma^2)$ , for  $i = 1, 2, \dots, I, j = 1, \dots, n_i$ . The test has a homoscedastic assumption, i.e., the error variance remains constant across all  $I$  groups.

Young and Bowman used a kernel-based smoothing method to approximate  $g_i$ . Assuming that  $h_i$  is the bandwidth for the  $i$ th regression function, they proposed to estimate  $g_i$  with

$$\hat{g}_i(x) = \frac{\sum_{j=1}^{n_i} K((x - x_{ij})/h_i) y_{ij}}{\sum_{j=1}^{n_i} K((x - x_{ij})/h_i)}, \quad (1)$$

which is sometimes referred to as the Nadaraya-Watson estimator of  $g_i$ .

Under the null hypothesis, one could obtain a common regression function by combining data from all groups

$$\hat{g}(x) = \frac{\sum_{i=1}^I \sum_{j=1}^{n_i} K((x - x_{ij})/h) y_{ij}}{\sum_{i=1}^I \sum_{j=1}^{n_i} K((x - x_{ij})/h)}, \quad (2)$$

where  $h$  is the common bandwidth for estimating  $g$ .

The test statistic that Young and Bowman proposed is analogous to the one-way ANOVA,

$$T_1 = \frac{\sum_{i=1}^I \sum_{j=1}^{n_i} [\hat{g}(x_{ij}) - \hat{g}_i(x_{ij})]^2}{\hat{\sigma}^2}, \quad (3)$$

where  $\hat{g}_i$  and  $\hat{g}$  are the group-specific and common curve estimators, and  $\hat{\sigma}^2$  is the pooled variance. To estimate  $\sigma^2$ , one uses  $\hat{\sigma}^2 = \frac{1}{N-I} \sum_{i=1}^I (n_i - 1) \hat{\sigma}_i^2$ , where  $N = \sum_{i=1}^I n_i$ . Similarly, the group-specific variance is estimated as

$$\hat{\sigma}_i^2 = \frac{1}{2(n_i - 1)} \sum_{j=1}^{n_i - 1} (y_{i, [j+1]} - y_{i, [j]})^2.$$

This test has been extended to comparisons of surface functions [17].

The construction of the test is intuitive, and its implementation straight-forward. The equal variance assumption, however, can be overly restrictive in some analytical situations. Additionally, when the explanatory variable  $x_j$  takes different values in the comparison groups, the power of the test often drops precipitously because the biases can no longer be canceled out under  $H_0$ ; see Tables 1–3 in Young and Bowman (1995)[16].

More recently, Park and colleagues considered a similar ANOVA type test statistic for multiple  $x$  values with a given bandwidth [18]. They obtained an empirical distribution for the maximum of the pointwise test statistics for controlling multiplicity. A visualization tool has been developed to show differences between curves at multiple locations. But simulation studies suggest that the test has type I error rates far below the nominal level, and very low power; see Table 1 of Park et al [18].

## 2.2 Kernel-based nonparametric methods

Dette and Neumeyer (2001) proposed another set of tests, all based on kernel smoothing techniques [19]. Expressing the curves as  $Y_{ij} = g_i(x_{ij}) + \epsilon_{ij}(x_{ij})$ , where  $i = 1, 2, \dots, I, j = 1, \dots, n_i$ , the authors introduced heteroscedastic errors  $\epsilon_{ij}(x_{ij}) \sim N(0, \sigma_i^2)$  into the model. The main hypothesis remains the same,  $H_0 : g_1 = g_2 = \dots = g_I$  vs  $H_1 : g_i \neq g_j$  for some  $i, j \in \{1, \dots, I\}$ .

The tests are subject to the following conditions: (1) The variances  $\sigma_i(\cdot)$  are continuous functions; (2) the design points  $x_{ij}$  satisfy  $\int_0^{x_{ij}} r_i(x) dx = \frac{j}{n_i}$  for a density function  $r_i$ , where  $j = 1, \dots, n_i$  and  $i = 1, \dots, I$ ; (3) the regression functions  $g_i(\cdot)$  are sufficiently smooth, i.e., 2 times continuously differentiable in the supporting space. And the Nadaraya-Watson estimators  $\hat{g}_i$  and  $\hat{g}$  are as previously defined.

One test ( $T_2$ ) compares the group-specific error variances against that of the combined sample, in a way that is analogous to one-way ANOVA

$$T_2 = \hat{\sigma}^2 - \frac{1}{N} \sum_{i=1}^I n_i \hat{\sigma}_i^2. \quad (4)$$

The second test ( $T_3$ ) directly assesses the distance between the group-specific curves and a common curve at  $x_{ij}$ , assuming that  $x_{ij}$  remain exactly the same across the groups,

$$T_3 = \frac{1}{N} \sum_{i=1}^I \sum_{j=1}^{n_i} [\hat{g}(x_{ij}) - \hat{g}_i(x_{ij})]^2. \quad (5)$$

The third test ( $T_4$ ) summarizes all pairwise  $L_2$ -distances of the estimated individual group curves,

$$T_4 = \sum_{i=1}^I \sum_{j=1}^{i-1} \int [\hat{g}_i(x) - \hat{g}_j(x)]^2 w_{ij}(x) dx, \quad (6)$$

where  $w_{ij}(\cdot)$  are positive weight functions. Asymptotic normality of the test statistics has been established under the null and fixed alternatives.

The above tests have been extended to comparison of two regression curves with different design points and heteroscedastic variances [20]. For comparison of two curves, the authors proved by using the empirical process theory that the two marked empirical processes converged to a centered Gaussian process at a rate of  $N^{-1/2}$  under the null; while under the alternative, the means of the two processes do not converge to zero. Hence, tests could be constructed based on either functions of the integrated squared residuals or the supremum of the absolute residuals of these two processes.

A practically important extension of Neumeyer and Dette's methods is the comparison of multiple curve functions [21]. These tests have also been extended to compare surface functions. Since the rates of convergence of the test statistics are slower [22], and the wild bootstrap procedure is consistent under more relaxed conditions, analysts often calculate  $p$ -values from the distribution of the test statistic under  $H_0$  using a wild bootstrap procedure [23].

### 2.3. Additive model-based tests

Zhang and Lin (2000) [24] considered testing the equivalence of two nonparametric functions in an additive mixed model for longitudinal data

$$Y_{ijk} = g_i(x_{ijk}) + \mathbf{s}_{ijk}^T \boldsymbol{\alpha}_i + \mathbf{Z}_{ijk}^T \mathbf{b}_{ij} + \epsilon_{ijk}, \quad (7)$$

where  $Y_{ijk}$  is the response from the  $j$ th subject in the  $i$ th group at the  $k$ th assessment, and  $\boldsymbol{\alpha}_i$  is a  $p \times 1$  vector associated with covariates  $s_{ijk}$ .

To test hypothesis  $H_0 : g_1 = g_2$  vs.  $H_1 : g_1 \neq g_2$ , the authors suggested the following test statistic

$$G\{\hat{g}_1(x), \hat{g}_2(x)\} = \int_{T_1}^{T_2} \{[\hat{g}_1(x) - \hat{g}_2(x)]^2\} dx, \quad (8)$$

where  $\hat{g}_1$  and  $\hat{g}_2$  are estimated by maximizing the penalized log-likelihood function.

The penalized likelihood under the semiparametric additive mixed model for an individual group is  $l(g_i, \boldsymbol{\alpha}_i; \mathbf{y}) - \frac{\lambda_i}{2} \int [g_i'(x)]^2 dx$ , where  $\lambda_i$  is the parameter that controls the smoothness of function  $g_i(x)$ . Note that the test statistic (8) is also an  $L_2$ -based distance measure, as that in (6). Expressing  $G$  in Equation (8) as a quadratic function of  $\mathbf{y}$ , Zhang and Lin approximated the distribution of  $G\{\hat{g}_1(x), \hat{g}_2(x)\}$  with a scaled  $\chi^2$  distribution using the moment matching technique.

We summarize the tests reviewed in Table 1, which highlights the key features of each method.

## 3. Curve comparison in semiparametric regression

In a low-dimensional regression analysis, comparing nonlinear effects amounts to a comparison of nonlinear curves and surfaces. As described in Table 1, the existing methods are often overly restrictive in their accommodation of heterogeneity and design points. For example, when we compare two nonlinear functions  $g_1(x)$  and  $g_2(x)$ , it is rather unrealistic to expect the two functions be evaluated at the exact same  $x$  values. Furthermore, if the functions are indeed different, it would not be reasonable to expect the two functions to have the same variance. These are the features that the existing methods and analytical software have not accommodated adequately.

To remedy, we propose a testing procedure in the usual context of semiparametric regression. For the convenience of discussion, we consider a comparison of curve functions  $g_i(x_{ij})$ , where  $x_{ij}$  denotes the value of independent variable of the  $j$ th subjects in the  $i$ th group. We are interested in testing hypothesis  $g_1 = g_2 = \dots = g_I$ . We note that the test can be extended to compare higher dimensional functions, although visualization may be difficult.

In order to test the above hypothesis, one needs to estimate the functions,  $g_1, g_2, \dots, g_I$ , as well as  $g$ . Most of the existing methods are developed based on kernel estimates. In practice, however, many analysts prefer various forms of regression splines, with or without smoothness penalty. In this paper, all theoretical results are derived under B-spline estimates. For all practical purposes, choices of spline basis functions are often less consequential. Here for curve comparison, we write the model as  $Y_{ij} = g_i(x_{ij}) + \epsilon_{ij}$ , where  $\epsilon_{ij} \sim N(0, \sigma_i^2)$ , and we write the group-specific function  $g_i$  as a B-spline function, i.e.,  $g_i(x) = \sum_{k=1}^{K_n+m} \gamma_k B_k^m(x)$ , where  $K_n$  is the number of internal knots with  $K_n = O(n^\nu)$ ,  $m$  is the order of the B-spline with  $m \geq 1$ ,  $\{\gamma_k\}_{k=1}^{K_n+m}$  is the set of B-spline coefficients, and  $\{B_k^m(x) : x \in [a, b]\}$  are B-spline basis functions. For higher dimensional functions  $g_i(x_1, x_2, \dots, x_d)$ , one could use tensor products, radial, or thin-plate splines [25][26][27].

### 3.1. Test statistics and comparison procedures

An intuitive way to compare two functions is to measure the distance between them. The  $L_2$  norm is a commonly used distance measure. As described previously, both Zhang and Lin (2000) and Dette and Neumeyer (2003) used the  $L_2$  norm in the construction of their test statistics. Herein, we reexamine the test statistic

$$T_{spline} = \frac{1}{N} \sum_{1 \leq i < m \leq I} \sum_{j=1}^{n_i} [\hat{g}_i(\mathbf{x}_{ij}) - \hat{g}_m(\mathbf{x}_{ij})]^2,$$

under B-spline estimates of  $\hat{g}_i$  and  $\hat{g}_m$  for testing the hypothesis  $H_0 : g_1 = g_2 = \dots = g_I$  vs  $H_1 : g_i \neq g_j$  for some  $i, j \in \{1, \dots, I\}$ . Theoretical properties of the test statistic are examined in Section 3.1.2.

In the absence of an asymptotic normal distribution, however, one has to devise a method to approximate the distribution of the test statistic under the null hypothesis. In the following section, we demonstrate how such an approximation can be done through a resampling procedure. Specifically, we show how to ascertain  $p$  values for the test statistic from a wild bootstrap procedure.

**3.1.1 A wild bootstrap-based comparison method**—For the standard linear regression models, it is usually sufficient to draw bootstrap samples from centralized residuals, because the errors are homoscedastic [28]. In the one dimensional case, the underlying model can be written as  $Y_{ij} = g_i(x_{ij}) + \epsilon_{ij}$ , where  $\epsilon_{ij} \sim N(0, \sigma_i^2)$ . We present  $g$  as a curve function here, although the method can be easily extended to higher dimensional situations.

To accommodate error heteroscedasticity, we consider a wild bootstrap procedure, which assures that the bootstrap error terms possess properties that are similar to those of the actual errors [29]. Another alternative approach is to use pairs bootstrapping, in which the analyst directly resamples from the joint empirical distribution function of  $\mathbf{Y}_j$  and  $\mathbf{x}_j$ , which are the vectors of the response and independent variables respectively. The computational burden of pairs bootstrap, however, tends to be greater especially if the dimension of  $\mathbf{x}_j$  is high [28].

Wild bootstrap has been used to resample the residuals of nonparametric regression models, as suggested by Härdle and Mammen (1993) [30], and Mammen (1993) [31]. The essence of wild bootstrap is to express the regression function as a conditional expectation of the observed response variable, i.e.  $E(Y_i^* | X_i = x_i) = g(x_i)$ , where  $Y_i^*$  is the bootstrap data. Since this method uses a single residual  $\hat{\epsilon}_i$  to estimate the conditional distribution  $I(Y_i - g(x_i) | X_i = x_i)$  of an arbitrary distribution ( $\hat{F}_i$  in the following), it is often referred to as the wild bootstrap.

Let  $V_i$  be a random variable following a two-point distribution  $\hat{F}_i$  such that  $E_{\hat{F}_i}(V_i) = 0$ ,  $E_{\hat{F}_i}(V_i^2) = 1$ , and  $E_{\hat{F}_i}(V_i^3) = 1$ . We define random independent quantities  $\epsilon_i^* = V_i \hat{\epsilon}_i \sim \hat{F}_i$ , and use  $(X_i, Y_i^* = \hat{g}(x_i) + \epsilon_i^*)$  as the bootstrap observations. We then create a new bootstrap test statistic  $T^*$ .

With the bootstrap samples, for a test at level  $\alpha$ , the null hypothesis will be rejected if  $T$  is greater than the corresponding quantile of the bootstrap distribution of the test statistic  $T^*$ , i.e.  $T > T_{[B(1-\alpha)]}^*$ , where  $T_{[B(1-\alpha)]}^*$  is the  $k$ th order value of the bootstrap statistic  $T^*$ . Härdle and Mammen (1993)[30] showed that under the null hypothesis,  $T^*$  estimated the distribution of  $T$  consistently, since the regression function with bootstrap data  $g^*(\cdot)$  had mean  $g(\cdot)$  for nonlinear models under the standard regularity conditions.

Using the wild bootstrap method, we propose the following procedure:

Step 1: Estimate function  $g_i(\mathbf{x})$  with  $\hat{g}_i(\mathbf{x})$ ,  $i = 1, 2, \dots, I$ , and compute the test statistic

$$T_{spline} = \frac{1}{N} \sum_{1 \leq i < m \leq I} \sum_{j=1}^{n_i} [\hat{g}_i(\mathbf{x}_{ij}) - \hat{g}_m(\mathbf{x}_{ij})]^2.$$

Step 2: Estimate the common function  $\hat{g}(\mathbf{x})$  from the combined sample and calculate the residuals  $\hat{\epsilon}_{ij} = y_{ij} - \hat{g}(\mathbf{x}_{ij})$ .

Step 3: For each  $\mathbf{x}_{ij}$ , draw a bootstrap residual  $\epsilon_{ij}^{(b)}$ , for  $b = 1, 2, \dots, B$ , where  $B$  is the number of bootstrap samples, from the two-point distribution with probability mass points  $\frac{1-\sqrt{5}}{2}\hat{\epsilon}_{ij}$  and  $\frac{1+\sqrt{5}}{2}\hat{\epsilon}_{ij}$ , occurring with probabilities  $\frac{5+\sqrt{5}}{10}$  and  $\frac{5-\sqrt{5}}{10}$  respectively, so that  $E(\epsilon_{ij}^{(b)}) = 0$ ,  $E(\epsilon_{ij}^{(b)2}) = \hat{\epsilon}_{ij}^2$  and  $E(\epsilon_{ij}^{(b)3}) = \hat{\epsilon}_{ij}^3$ .

Step 4: Generate a bootstrap sample  $(x_{ij}, Y_{ij}^{(b)})$  from  $Y_{ij}^{(b)} = \hat{g}(\mathbf{x}_{ij}) + \epsilon_{ij}^{(b)}$ .

Step 5: From this sample, estimate the  $b$ th bootstrap regression function  $\hat{g}_i^{(b)}$ , and calculate the test statistic  $T^{(b)}$  as in the original  $T_{spline}$  calculation.

Step 6: Repeat Step 3 to 5  $B$  times, and use the  $B$  generated values of the test statistics,  $T_{spline}^* = (T^{(1)}, T^{(2)}, \dots, T^{(B)})$ , to determine the quantiles of the distribution of the test statistic. For a test at significance level  $\alpha$ , the null hypothesis is rejected if  $T_{spline}$  is greater than the corresponding  $(1 - \alpha)$ th quantile of the bootstrap distribution of  $T_{spline}^*$ .

**3.1.2. Consistency of the test against fixed alternatives**—In this section we show that the proposed test is consistent against any fixed alternatives. We also provide the optimal number of internal knots for the B-spline that leads to the best possible rate of convergence. We first rewrite the model in a slightly more general form; here we use  $X$  instead of  $x$  to emphasize the random nature of the independent variable.

We write the true model as follows

$$Y_j = g_0(X_j) + \epsilon_j = g_0(X_j) + \sigma \epsilon_j,$$

where  $g_0$  is an unknown function of interest,  $(Y_j, X_j), j = 1, \dots, n$ , are i.i.d. random variables independent of the error term  $\epsilon_j \sim \mathcal{N}(0, 1)$ . For simplicity and without loss of generality, we assume that the covariate  $X_j \in \mathbb{X}$  a.s., where  $\mathbb{X} = [0, 1]$ . As defined,  $\mathbb{X}$  is a compact subset in  $\mathbb{R}$ .

As previously described, a B-spline estimate of  $g(x) = \sum_{k=1}^{K_n+m} \gamma_k B_k^m(x)$  can be achieved by minimizing the objective function

$$\frac{1}{n} \sum_{j=1}^n [Y_j - g(X_j)]^2,$$

or, equivalently, by maximizing

$$\mathbb{M}_n(g) \equiv \mathbb{P}_n m_g = 2\mathbb{P}_n(g - g_0)e - \mathbb{P}_n(g - g_0)^2,$$

where  $\mathbb{P}_n m_g$  denotes the empirical process indexed by the function  $m_g$ , i.e.

$$\mathbb{P}_n m_g = \frac{1}{n} \sum_{j=1}^n m_g(X_j) = -\frac{1}{n} \sum_{j=1}^n [Y_j - g(X_j)]^2.$$

Direct maximization of the above objective function over the full infinite-dimensional parameter space  $\mathcal{G}$  may lead to inconsistent estimates[32]. There-fore, one has to use the sieve M-estimation framework by considering the spaces of B-spline functions



$$\mathcal{S}_n(D_n, K_n, m) = \left\{ g_n : g_n(x) = \sum_{k=1}^{K_n+m} \gamma_k B_k^m(x) \in S_n(D_n, K_n, m), x \in [0, 1] \right\},$$

where  $D_n = \{d_1, \dots, d_{K_n}\}$  is a set of partition points for the set  $[0, 1]$ ,  $K_n$  is the number of internal knots with  $K_n = O(n^\nu)$ ,  $m$  is the order of the B-spline with  $m \geq 1$ ,  $\{\gamma_k\}_{k=1}^{K_n+m}$  is the set of the unknown coefficients or control points for the B-spline,  $\{B_k^m(x) : x \in [a, b]\}$  are the basis functions, and  $S_n(D_n, K_n, m)$  is the space of polynomial splines on a partition  $D_n$  with  $K_n$  internal knots and of order  $m$ . Then, the sieve estimator  $\hat{g}_n$  of  $g_0$  satisfies

$$\mathbb{M}_n(\hat{g}_n) \geq \mathbb{M}_n(g) \text{ for all } g \in \mathcal{S}_n,$$

that is  $\hat{g}_n$  maximizes  $g \mapsto \mathbb{M}_n(g)$  over the sieve space  $\mathcal{S}_n(D_n, K_n, m)$ .

For simplicity of presentation, we consider the special case of the two-sample comparison. We assume the following regularity conditions:

- C1. The error  $e_j$  has a distribution with zero mean and sub-exponential tails, i.e., tails bound by the supremum of an empirical process. Also,  $e$  and the covariate  $X$  are independent.
- C2. The parameter space  $\mathcal{S}_i \ni g_{0,i}$ ,  $i = 1, 2$ , contains functions uniformly bounded by  $C^{-1/2}$  on  $[0, 1]$ , with bounded  $p$ th derivatives, for fixed  $p \geq 1$ , with the first derivative being continuous.
- C3. The number of internal knots satisfies  $K_n = O(n^\nu)$ , such that

$$\max_{1 \leq k \leq K_n+1} \{d_k - d_{k-1}\} = O(n^{-\nu}).$$

- C4. The sample sizes of the two groups satisfy

$$\frac{n_1}{n_1 + n_2} \rightarrow \lambda \in (0, 1),$$

as  $\min(n_1, n_2) \rightarrow \infty$ .

**Theorem 1.** Assuming conditions C1-C4 hold, we consider a wild-bootstrap procedure that of level  $\alpha$  asymptotically. Then, the proposed test is consistent against any fixed alternative hypothesis: If  $\pi_n(\theta_\delta)$  is the power function of the test under the fixed alternative hypothesis  $\theta_\delta$ , then  $\pi_n(\theta_\delta) \rightarrow 1$  as  $n \rightarrow \infty$ .

A sketch of the proof is provided in Appendix A. Even though the proposed test is, by Theorem 1, consistent against any fixed alternative hypothesis, the asymptotic distribution of the test statistic is difficult to derive. The difficulty stems from the fact that the convergence rate of the B-spline estimator of  $g_1$  and  $g_2$  is

$$d(\hat{g}_{n_i, i}, g_{0, i}) = O_p\left(\frac{p}{n^{1+2p}}\right), \quad i = 1, 2,$$

where  $d(g_1, g_2) = \{E[g_1(X) - g_2(X)]^2\}^{1/2}$ . Note that the above rate is the optimal convergence rate for nonparametric regression and is achieved if one sets  $v = 1/(1 + 2p)$ . This convergence rate is slower than the usual  $\sqrt{n}$  rate for parametric models, even though it is the optimal rate in nonparametric regression.

In the absence of an analytically derived asymptotic distribution, we assessed the performance of the wild-bootstrap procedure through extensive simulations. Results of the simulation experiments are presented in Section 5.

**3.1.3. Extending the testing procedure to correlated data**—The same testing procedures can be modified and extended to the analysis of correlated data. A key requirement for the modification is the preservation of the correlation structure that exists within each subject. The following algorithm is a natural extension and it performs a Cholesky decomposition on the estimated covariance structure [33]. Combining with the estimated regression functions, we generate the bootstrap sample  $Y_{ijk}^{(b)}$ .

The algorithm is described below.

Step 1: Obtain group-specific estimates  $\hat{g}_i(\mathbf{x})$  of function  $g_i$  by using semiparametric mixed effect models, and then compute the test statistic  $T_{splcorr}$ .

Step 2: Obtain a common regression function estimate  $\hat{g}(\mathbf{x})$  for the combined sample.

Step 3: For  $i = 1, \dots, I$ , obtain the group-specific covariance matrix estimate  $\hat{\mathbf{R}}_i$  from the fitted group-specific model and ascertain the residuals  $\hat{\eta}_i = (\hat{\eta}_{i11}, \hat{\eta}_{i12}, \dots, \hat{\eta}_{i n_i})$ .

Step 4: Perform a Cholesky decomposition on  $\hat{\mathbf{R}}_i$  so that  $\hat{\mathbf{R}}_i = \hat{\mathbf{L}}_i \hat{\mathbf{L}}_i^T$ , where  $\hat{\mathbf{L}}_i$  is a lower triangular matrix. We then obtain

$$\hat{\mathbf{e}}_i = (\hat{e}_{i11}, \hat{e}_{i12}, \dots, \hat{e}_{i n_i}) = \hat{\mathbf{L}}_i^{-1} \hat{\eta}_i;$$

and calculate the “whitened” residuals  $\tilde{\mathbf{e}}_i = \hat{\mathbf{e}}_i - \frac{1}{n_i} \sum_{j=1}^{n_i} \hat{\mathbf{e}}_i$ .

Step 5: Draw a random sample from the “whitened” residuals  $\tilde{\mathbf{e}}_i$  and calculate

$\hat{\eta}_i^{(b)} = \hat{\mathbf{L}}_i \tilde{\mathbf{e}}_i$ . We use  $Y_{ijk}^{(b)} = \hat{g}(\mathbf{x}_{ijk}) + \hat{\eta}_{ijk}^{(b)}$  as the bootstrap sample.

Step 6: Calculate the test statistic  $T^{(b)}$  with the bootstrap sample. Repeat Steps 3 to 5  $B$  times to construct an empirical distribution of the statistic with bootstrap samples, and then calculate the  $p$  value from the tail area of the empirical distribution.

## 4. Software development

We published an R package `gamm4.test` in CRAN to make the proposed testing procedures available to practitioners. The two main functions are `gam.grptest` for comparisons of nonlinear functions with cross-sectional data, and `gamm4.grptest` for comparisons involving correlated data. Key features of this package are:

- a. It utilizes a syntax that is consistent with `mgcv` and `gamm4`, two packages that are often used for fitting semiparametric regression models. Users familiar with those packages can perform comparisons with `gam.grptest` and `gamm4.grptest`.
- b. The package performs parallel computing with an automatic detection of the numbers of available CPU cores for enhanced computational efficiency.
- c. The R package includes a data visualization function `plot.gamtest` that allows users to visually examine the fitted curves and surfaces. The graphics are produced by the R package `plotly`. With option `type = plotly.persp`, the users can create 3-dimensional interactive plots. Finally, setting the option `test.statistic = TRUE` generates the empirical distribution of the test statistic under the null hypothesis of equal regression functions.

Computational efficiency of `gamm4.test` in the analysis of the example data in the package was assessed on a computer with Intel(R) Core(TM) i5-3470, CPU @3.20GHz, 64-bit operating system, and 4 CPU cores. The computing time is summarized in Table 2.

To enhance the usability of the testing methods, we also created an interactive R Shiny interface for the `gamm4.test` package. This interface allows analysts that do not use R to access the testing procedure through a web link. See <https://heather.shinyapps.io/shinygamm4/> for the app. See <https://youtu.be/SHqaZXSLaMw> for a related youtube tutorial.

## 5. Simulation Studies

We conducted a series of simulation studies to verify the theoretical results and to examine the operating characteristics of the proposed tests, in comparison with the other  $L_2$ -based methods. We additionally investigated the influences of the number of knots on the tests.

### 5.1. Curve comparisons

For curve comparisons, we considered the following model

$$Y_{ij} = g_{id}(x_{ij}) + \epsilon_{ij}, \quad (9)$$

where  $i = 1, 2; j = 1, \dots, n_i$ .

In this simulation, we generated values of the independent variable  $x_{ij}$  from  $Unif[0, 1]$ , with sample sizes  $n_1$  and  $n_2$  for the two comparison groups. The nonlinear functions for the two groups were specified as  $g_1(x) = 2x \exp(2-4x) - 2x + 0.5$  and  $g_2(x) = 2x \exp(2-4x) - 0.5$ , with  $g(x) = (4x \exp(2-4x) - 2x)/2$ . More generally, we considered  $g_{id}(x_{ij}) = (d/10)g_f(x_{ij}) + (1 - d/10)g(x_{ij})$ , with  $d = 0, 1, 2, 3$ , where  $d$  controlled the distance between the two group-specific

functions. For example,  $d=0$  corresponded to the situation where the two groups shared the same regression function; as  $d$  increased, the functions grew further apart. These functions were plotted in Appendix B; see Supplemental Figure 1.

Values of the dependent variables  $Y_{ij}$  were generated from Equation (9) with standard errors  $\sigma_1$  and  $\sigma_2$ , i.e.  $\epsilon_{1j} \sim N(0, \sigma_1^2)$ ,  $\epsilon_{2j} \sim N(0, \sigma_2^2)$ . Comparisons of the nonlinear functions were carried out under the following three conditions: (1) a distance parameter  $d=0, 1, 2, 3$ ; (2) sample sizes  $(n_1, n_2)=(125, 125), (216, 216)$ , and  $(512, 512)$ ; (3) values of the error standard deviations  $(\sigma_1, \sigma_2)=(0.20, 0.15)$  and  $(0.25, 0.20)$ .

With the generated data sets, we comparatively evaluated the performance of five discussed methods:

- Method 1: The proposed testing method with cubic B-spline regression bases for curve estimation, and a wild bootstrap procedure for p-value calculation.
- Method 2: The proposed testing method with penalized cubic spline bases for curve estimation, with a wild bootstrap procedure for p-value calculation by using the gam function in the R package mgcv. Numbers of knots were set to the default value, which was determined by a generalized cross-validation (GCV) method.
- Method 3: Kernel smoothing based on the  $L_2$  distance test statistic, followed by a wild bootstrap procedure [19].
- Method 4: The testing method based on variance estimator [19].
- Method 5: Young and Bowman's (1995) method [16], which calculates the  $p$  value by matching with a scaled chi-square distribution.

For each simulation setting, we generated a total of 1,000 testing datasets. For each dataset, we performed the proposed test on 200 wild bootstrap samples to calculate the  $p$  value. We calculated the rate of rejection in the 1,000 simulated samples at a significance level of 0.05.

Rejection rates for curve comparison under the null and alternative hypotheses are reported in Appendix B; see supplemental Table 1. When  $d=0$ , the rejection rates are Type I error rates; when  $d=1, 2, 3$ , the rejection rates represent the power of the tests. In comparison with other testing methods, the proposed tests in general had an excellent control of Type I error rates. As  $d$  increased, the power of rejecting the null hypothesis increased as well. The power of the new tests was comparable to, if not slightly better than, that of the existing tests. We noticed that Method 5 showed a slightly higher type I error rates than others. On the other hand, the Method 4 exhibited a tighter type I error control, while having considerably lower power. For the proposed testing methods, B-splines and penalized splines produced similar results. As previously discussed, we set the number of knots to  $\sqrt[3]{n_1}$ . In the simulation studies we observed that when an unpenalized semiparametric model was used, incorrect number of knots selection could lead to substantial bias and hence inflated the Type I error rates. The penalized semiparametric regression estimates, however, were generally very robust.

To verify the consistency theory in Section 3, we further examined the rejection rates (i.e., power) of the test as the sample size increased, under a fixed alternative hypothesis. We showed that when the distance between the null and alternative hypotheses was set to  $d = 1$ , the power increased with the sample size. The power approached to 1 when  $n_1 = n_2 = 1000$ . See Figure 2 in Appendix B.

Besides of the five methods listed above, we also compared the proposed methods with the test proposed by Kulasekera (1995) [13]. The latter's performance has left much to be desired: In the tested settings, the type I error rates were close to zero, and the power was low as well. The suboptimal performance could be due to the simulation settings we chose, where the functional curves were close and data variability was large. The proposed tests, on the other hand, performed well in such situations. We omitted the results of Kulasekera's test from the summary table.

We further examined the performance of the tests in situations of three group comparison. The nonlinear functions of the three groups were specified respectively as  $g_1(x) = 2x \exp(2-4x) - x - (d_1/10)(x-0.5)$ ,  $g_2(x) = 2x \exp(2-4x) - x + (d_2/10)(x-0.5)$  and  $g_3(x) = 2x \exp(2-4x) - x$ , with  $d_j = 0, 1, 2, 3$  and  $i = 1, 2$ , where  $d_1, d_2$  respectively controlled the distances between  $g_1$  and  $g_3$ , and  $g_2$  and  $g_3$ . Comparisons of the nonlinear functions were carried out with various distance values of  $d_1$  and  $d_2$ , as well as sample sizes. For each setting, we generated a total of 500 testing datasets.

See Supplemental Table 2 in Appendix B for rates of rejection in 500 simulated samples at a significance level of 0.05. Similar to two curve comparisons, when  $d_1 = d_2 = 0$ , the proposed tests in general had a good control of type I error rates. As  $d_j$  increased, regardless the equal or unequal distances between the curves, the power of rejecting the null hypothesis increased with the sample size.

## 5.2. Surface comparisons

Simulations for surface comparisons were carried out in a similar manner. We used the following surface functions in the simulation:

- a.  $g_1(\mathbf{x}) = g_2(\mathbf{x}) = \sin(2\pi x_1) + \cos(2\pi x_2)$
- b.  $g_1(\mathbf{x}) = g_2(\mathbf{x}) = 2x_1^2 + 3x_2^2$
- c.  $g_1(\mathbf{x}) = g_2(\mathbf{x}) = \exp(-x_1^2 - x_2^2)$
- d.  $g_1(\mathbf{x}) = \sin(2\pi x_1) + \cos(2\pi x_2)$   $g_2(\mathbf{x}) = \sin(2\pi x_1) + \cos(2\pi x_2) + x_1$
- e.  $g_1(\mathbf{x}) = 2x_1^2 + 3x_2^2$   $g_2(\mathbf{x}) = 2x_1^2 + 3x_2^2 + \sin(2\pi x_1)$
- f.  $g_1(\mathbf{x}) = \exp(-x_1^2 - x_2^2)$   $g_2(\mathbf{x}) = \exp(-x_1^2 - x_2^2) + \sin(2\pi x_1)$

Scenarios a-c represented situations under the null hypothesis, i.e., where the surfaces were the same; Scenarios d-f corresponded to various alternative hypotheses. The independent variables  $x_1$  and  $x_2$  were simulated from independent  $Unif[0, 1]$  with sample size  $n_1$  and  $n_2$  for each group. The dependent variables  $Y_{ij}$  were generated from the above functions with

standard errors  $\sigma_1$  and  $\sigma_2$ . i.e.  $Y_{ij} = g_i(x_{ij}) + \epsilon_{ij}$ , where  $i = 1, 2; j = 1, \dots, n_i$ ,  $\epsilon_{1j} \sim N(0, \sigma_1^2)$ ,  $\epsilon_{2j} \sim N(0, \sigma_2^2)$ .

We conducted the simulation under the following parameter settings: (1) Three sample size settings of  $(n_1, n_2)$  as (125, 125), and (216, 216), (512, 512); (2) Two different values of standard errors  $(\sigma_1, \sigma_2)$  for each function. For each simulation setting, we generated 500 datasets. For each dataset, we tested the new method with 300 wild bootstrap resamples. We calculated the rejection rate based on the 500 simulated datasets at  $\alpha = 0.05$ .

Type I error rates for function pairs of a-c are reported in Appendix B; see supplemental Table 4; powers for function pairs in d-f are reported in Table 5 of Appendix B. The Type I error rates of the proposed tests were generally good and power was superior than the existing tests. Similar to the simulation studies for curve comparisons, the penalized semiparametric model with the default number of knots from 'GCV' method showed a performance similar to the tests using semiparametric estimating methods and  $\sqrt[3]{n_i}$  number of knots. Numbers of knots had relatively minor influences on the testing performance.

In comparison with the existing methods, we found that in general the proposed methods had reasonable Type I error control as expected. The power was either comparable to or superior than that of the other methods.

### 5.3. Tests with correlated data

We considered the following models for correlated data

$$Y_{ijk} = g_{id}(x_{ijk}) + b_{ij} + \epsilon_{ijk},$$

where  $i = 1, 2; j = 1, \dots, n_i; k = 1, 2, 3$ . As previously presented in Equation (9), we used  $Y_{ijk}$  to indicate the measure on the  $k$ th occasion in subject  $j$  from group  $i$ . Values of the independent variable  $x_{ijk}$  were generated from independent  $Uni[0, 1]$ . Values of the dependent variable  $Y_{ij}$  were generated based on the above functions with random effect  $b_i \sim N(0, \sigma_i^2)$  and the i.i.d. random error  $\epsilon_{ijk} \sim N(0, \sigma_j)$ . We used the same regression functions in Section 5.1, where the two curves gradually grew apart with an increasing  $d$  (see Appendix Figure 1).

We considered the following parameter settings: (1)  $d=0, 1, 2$ ; (2) three sample size settings  $(n_1, n_2) = (50, 60), (100, 120),$  and  $(150, 160)$  and all with three repeated measures; (3) three different combinations of the standard deviations of the random intercept and the i.i.d random variable as  $(\sigma'_1, \sigma'_2, \sigma_1, \sigma_2) = (0.2, 0.15, 0.04, 0.05), (0.2, 0.15, 0.10, 0.12),$  and  $(0.25, 0.20, 0.10, 0.12)$ . We used penalized semiparametric mixed regression to estimate the curves.

The simulation results are reported in supplemental Table 6 of Appendix B. Results suggested a relatively tight Type I error rate control and good power. Zhang's (2000) scaled chi-square test is the only existing comparative test for correlated data [24]. We compared the performance of the new test with that of the scaled chi-square test. Under

sample sizes  $(n_1, n_2) = (50, 50)$  and  $(100, 100)$ , we performed the test using 1) the same  $x$  values for two groups; 2) slightly different  $x$  ( $x_2 = x_1 + Unif(0, 0.05)$ ); 3) completely random and independent  $x_1, x_2$  for two groups. Simulation was repeated for 200 times under each scenario and the results were shown in Appendix B Table 7.

The proposed test clearly outperformed the scaled  $\chi^2$  test. When the two groups had the same values in  $x$ , the scaled  $\chi^2$  test had Type I error rates that were lower than the nominal level. The power was generally lower as well. The scaled chi-square test was not designed for situations of randomly distributed independent variables so the power deteriorated when we introduced different  $x$  values between the two groups.

## 6. Real data applications

To illustrate the proposed testing procedures, we analyzed two data sets from a large observational study. The original study was designed to examine the factors related to blood pressure development in children. Detailed study protocols were published elsewhere [34] [35]. Briefly, healthy children between 5 and 17 years of age were recruited from schools in Indianapolis, Indiana. Blood pressure, height, and weight were measured twice a year from the study participants. Blood and overnight urine samples were collected. The study protocol was approved by a local Internal Review Board. Informed consent was obtained from study participants, or their parents when appropriate.

### 6.1. Comparisons of weight growth curves

We compared the weight growth curves between blacks and whites within each sex, and between boys and girls within each race. We write the model as

$$Weight_{ij} = g_i(Age_{ij}) + \epsilon_{ij},$$

where  $i$  indexes the groups and  $j$  the  $n_j$  observations within each group. Here we used the baseline assessment data to examine the weight-age relationships in the four sex and race combinations. The sizes of the four groups were: 205 black boys, 311 white boys, 232 black girls, and 289 white girls. We performed comparisons by testing the hypotheses  $H_0 : g_1 = g_2$  vs  $H_1 : g_1 \neq g_2$ , where  $g_1, g_2$  are the weight growth curves between the sexes within each race group, or weight growth functions between the races within each sex group.

We first estimated the weight growth curves of the groups as part of the preliminary analysis. See scatter plots in Figure 1(a). We reported the  $p$  values of the four competing testing methods in Table 3. We presented the curve estimates with 95% pointwise confidence intervals from the semiparametric regression model (Generalized Cross-Validation for selecting smoothing parameter and thin-plate penalized basis function) in Figure 1(b). The nonparametric smoothing curves by loess produced curve estimates that were similar to the semiparametric regression estimates.

Testing results from the semiparametric spline-based estimating method were consistent with the curve estimations shown in Figure 1. The tests showed that the weight-for-age curves were significantly different between white and black girls. In our sample, the black

girls gained more weight around ages 12 and 13 than their white peers, but the two curves converged gradually at age 14 years. The confidence intervals became wider after age 15, possibly due to the reduced sample sizes. Similar patterns were seen in the height-for-age curves.

For surface comparison, we considered weight as a function of age and height. We wrote the model as

$$Weight_{ij} = g_i(Age_{ij}, Height_{ij}) + \epsilon_{ij},$$

where  $i$  is the index for the sex-race group, and  $j$  is the index for a specific subject,  $j = 1, 2, \dots, n_i$  within the group. We compared simultaneous effects of height and age on weight, among the four race-sex groups. The p-values of the four types of tests were summarized in Table 3 and the corresponding contour plots were presented in Figure 2. No statistically significant differences were detected using the four tests.

## 6.2. Hormonal influences on blood pressure

Blood pressure is regulated by hormones in the renin-angiotensin-aldosterone system (RAAS). An essential product of RAAS is aldosterone, a mineralocorticoid hormone. Aldosterone acts on the epithelial sodium channel (ENaC) to help retain sodium. Increased sodium load causes extracellular fluid volume (ECFV) expansion, which in turn leads to blood pressure elevation. Recent biological experiments have shown that in the American population, blacks are more responsive to the stimulation of aldosterone in comparison to whites [35]. As a result, blacks tend to have greater levels of ECFV, which helps to suppress renin secretion. Renin, together with potassium, helps production of aldosterone [36]. This process is essential for the maintenance of blood pressure [37].

In this analysis, we examined the simultaneous influences of plasma renin activity (PRA) and plasma aldosterone concentration (PAC) on systolic blood pressure (BP):

$$BP_{ij} = g_i(\log(PRA)_{ij}, \log(PAC)_{ij}) + \epsilon_{ij},$$

where  $i$  is the index for the race group, and  $j$  is the index for a specific subject,  $j = 1, 2, \dots, n_i$  within the group. We were interested in comparing the surface functions in blacks and in whites. A quick visual examination (Figure 3) showed that blood pressure in whites ( $n=313$ ) was on average lower than that in blacks ( $n=184$ ). In whites, PRA and PAC were not significantly correlated with blood pressure, as one would expect in a steady-state sample. But in blacks, lower PRA and higher PAC were associated with a higher systolic blood pressure, which suggested an increased blood pressure sensitivity to aldosterone. We compared the two surface functions and obtained a p-value of 0.054, based on 500 bootstrap resamples. The comparison generated a  $p$  value that was close to, but did not reach the commonly accepted threshold of 0.05. Other testing methods, including Young and Bowman's (1995) tests (Method 4 and 5 in Section 5.1) gave p-values of 0.15 and 0.73, while the  $L_2$  distance-based test using kernel smoothing (Method 3 in Section 5.1) provided a p-value of 0.03. Indication for the racial difference in this analysis is generally consistent



with the experimental evidence from drug-induced hyperaldosteronism in human subjects [35]; the actual power of the tests, however, is likely influenced by the sample size and data variability.

## 7. Discussion

Statistical analysis of biomedical data is never complete without a proper test. In analysis of lower-dimensional nonlinear functions, inference typically involves comparisons of curves and surfaces. In parametric analysis where the functions are fully specified, inference is generally straightforward and can be carried out in the usual likelihood-based framework. In nonparametric or semiparametric analyses, due to lack of knowledge of the true functional forms of the relationships, hypotheses cannot be formulated solely on prespecified parameters. Analysts, therefore, can no longer rely on likelihood-based tests. Standard software packages or functions typically do not produce comparison of interest. In practice, estimation and inference of the functional curves and surfaces are often done separately, in part due to the lack of integration of estimation and inference tools and common programming syntax. In the present paper, we propose new testing procedures based on the  $L_2$  distance. We show that the proposed tests are consistent against any fixed alternative hypothesis. To evaluate the level of statistical significance we provide a set of bootstrap testing methods. We have developed an R package to assist analysts who are interested in using the tests. For those who do not use R, we present an R Shiny interface to wrap around the software package so that analysts could directly upload their data to a web server and implement the tests through interactive web-based operations.

Extensive simulation studies show that, in comparison to the existing methods, the proposed tests have good control of type I error rate and excellent power. Despite our use of computer intensive methods such as the wild bootstrap, the procedures are generally quite efficient. Our own testing of the method with real data shows that the software package is easy to operate and it is flexible for accommodating covariates and repeated measures. We showed that the testing procedure possesses the property of consistency, a necessary condition for bootstrap to work. The rate of convergence, however, has not reached the optimum of  $n^{1/2}$ . The empirical evidence from our simulation has nonetheless supported the good performance in finite sample situations. Notwithstanding this limitation, we put forward a computational tool for the comparison of curves and surfaces in nonparametric or semiparametric analyses.

## Supplementary Material

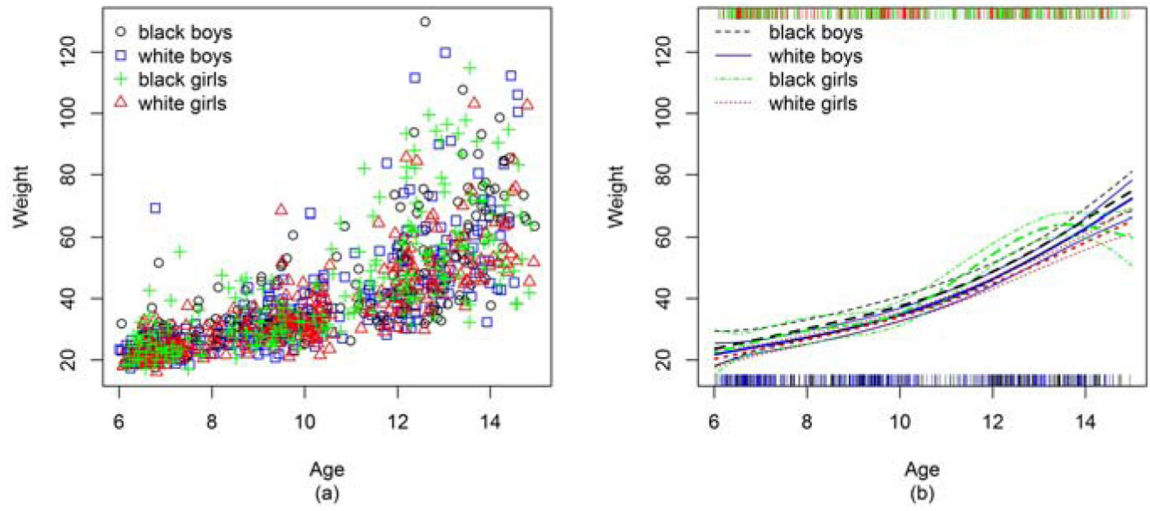
Refer to Web version on PubMed Central for supplementary material.

## References

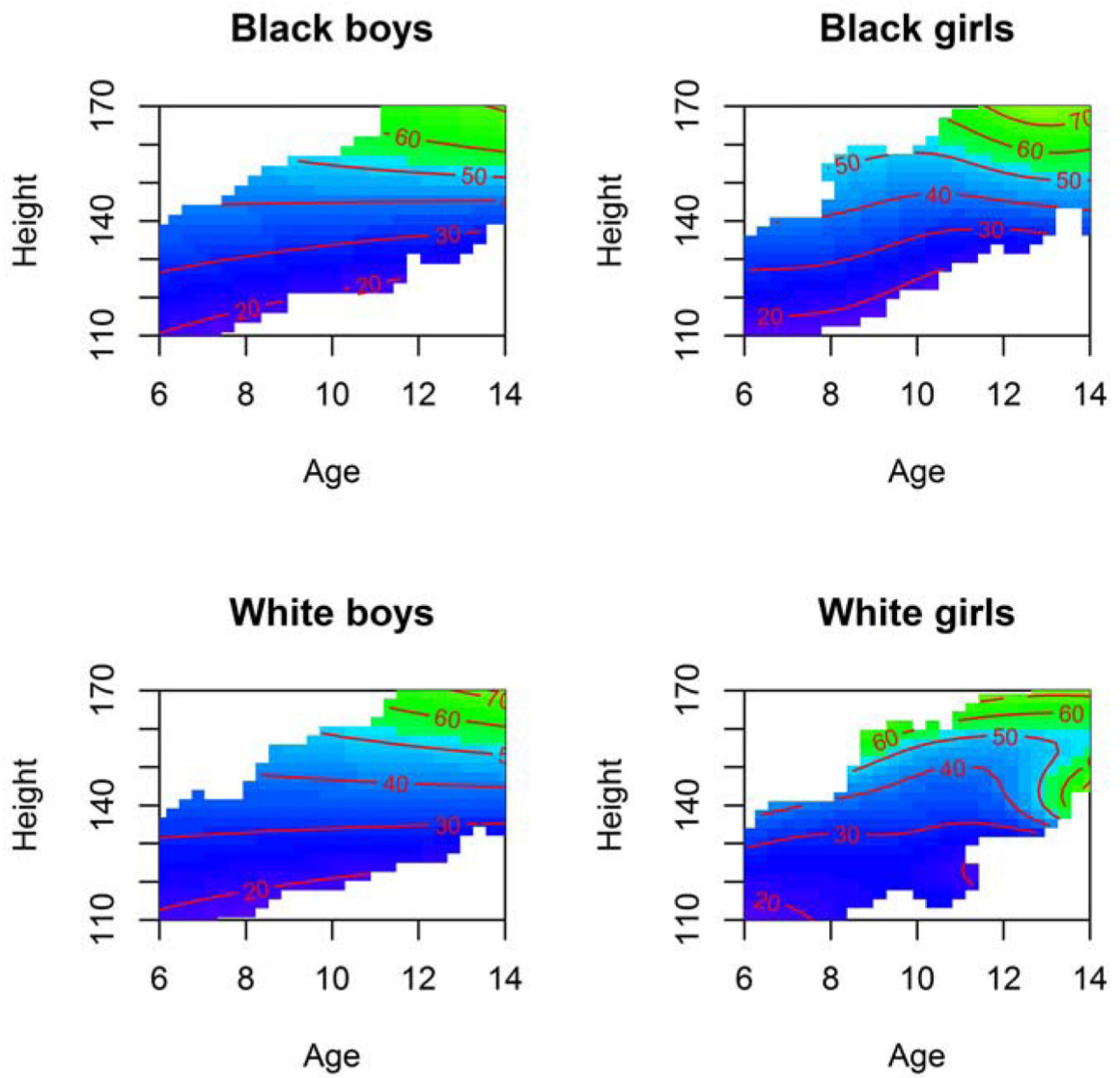
- [1]. Green P, Silverman B, Kernel nonparametric regression and generalized linear models: A roughness penalty approach (1994).
- [2]. Fan J, gijbels I, Local Polynomial Modelling and Its Applications.
- [3]. Ogden T, Essential wavelets for statistical applications and data analysis, Springer Science & Business Media, 2012.

- [4]. Wahba G, Spline models for observational data, Vol. 59, Siam, 1990.
- [5]. Gu C, Smoothing spline ANOVA models, Vol. 297, Springer Science & Business Media, 2013.
- [6]. Marx BD, Eilers PH, Generalized linear regression on sampled signals and curves: a p-spline approach, *Technometrics* 41 (1) (1999) 1–13.
- [7]. Eubank RL, Nonparametric regression and spline smoothing, CRC press, 1999.
- [8]. Bde Boor C, A practical guide to splines, revised edition (2001).
- [9]. Hastie T, Tibshirani R, Generalized additive models, London: Chapman and Hall (1990) 137–173.
- [10]. Delgado MA, Testing the equality of nonparametric regression curves, *Statistics & probability letters* 17 (3) (1993) 199–204.
- [11]. Härdle W, Marron JS, et al. , Semiparametric comparison of regression curves, *The Annals of Statistics* 18 (1) (1990) 63–89.
- [12]. Fan J, Lin S-K, Test of significance when data are curves, *Journal of the American Statistical Association* 93 (443) (1998) 1007–1021.
- [13]. Kulasekera K, Comparison of regression curves using quasi-residuals, *Journal of the American Statistical Association* 90 (431) (1995) 1085–1093.
- [14]. Eubank R, Li C, A diagnostic test for parallelism, *Journal of Statistical Sciences* (2008) 13–29.
- [15]. Dette H, Dhar SS, Wu W, Identifying shifts between two regression curves, arXiv preprint arXiv:1908.04328.
- [16]. Young SG, Bowman AW, Non-parametric analysis of covariance, *Biometrics* (1995) 920–931.
- [17]. Bowman AW, Comparing nonparametric surfaces, *Statistical Modelling* 6 (4) (2006) 279–299.
- [18]. Park C, Hannig J, Kang K-H, Nonparametric comparison of multiple regression curves in scale-space, *Journal of Computational and Graphical Statistics* 23 (3) (2014) 657–677.
- [19]. Dette H, Neumeier N, et al. , Nonparametric analysis of covariance, *the Annals of Statistics* 29 (5) (2001) 1361–1400.
- [20]. Neumeier N, Dette H, et al. , Nonparametric comparison of regression curves: an empirical process approach, *The Annals of Statistics* 31 (3) (2003) 880–920.
- [21]. Pardo-Fernández JC, Van Keilegom I, González-Manteiga W, Testing for the equality of k regression curves, *Statistica Sinica* (2007) 1115–1137.
- [22]. Hall P, Hart JD, Bootstrap test for difference between means in nonparametric regression, *Journal of the American Statistical Association* 85 (412) (1990) 1039–1049.
- [23]. Wang X-F, Ye D, On nonparametric comparison of images and regression surfaces, *Journal of statistical planning and inference* 140 (10) (2010) 2875–2884. [PubMed: 20543891]
- [24]. Zhang D, Lin X, Sowers M, Semiparametric regression for periodic longitudinal hormone data from multiple menstrual cycles, *Biometrics* 56 (1) (2000) 31–39. [PubMed: 10783774]
- [25]. Wood SN, Thin plate regression splines, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65 (1) (2003) 95–114.
- [26]. Wood SN, Low-rank scale-invariant tensor product smooths for generalized additive mixed models, *Biometrics* 62 (4) (2006) 1025–1036. [PubMed: 17156276]
- [27]. Liu H, Tu W, et al. , A semiparametric regression model for paired longitudinal outcomes with application in childhood blood pressure development, *The Annals of Applied Statistics* 6 (4) (2012) 1861–1882.
- [28]. Freedman DA, et al. , Bootstrapping regression models, *The Annals of Statistics* 9 (6) (1981) 1218–1228.
- [29]. Liu RY, et al. , Bootstrap procedures under some non-iid models, *The Annals of Statistics* 16 (4) (1988) 1696–1708.
- [30]. Härdle W, Mammen E, et al. , Comparing nonparametric versus parametric regression fits, *The Annals of Statistics* 21 (4) (1993) 1926–1947.
- [31]. Mammen E, et al. , Bootstrap and wild bootstrap for high dimensional linear models, *The Annals of Statistics* 21 (1) (1993) 255–285.
- [32]. Shen X, Wong WH, Convergence rate of sieve estimates, *The Annals of Statistics* (1994) 580–615.

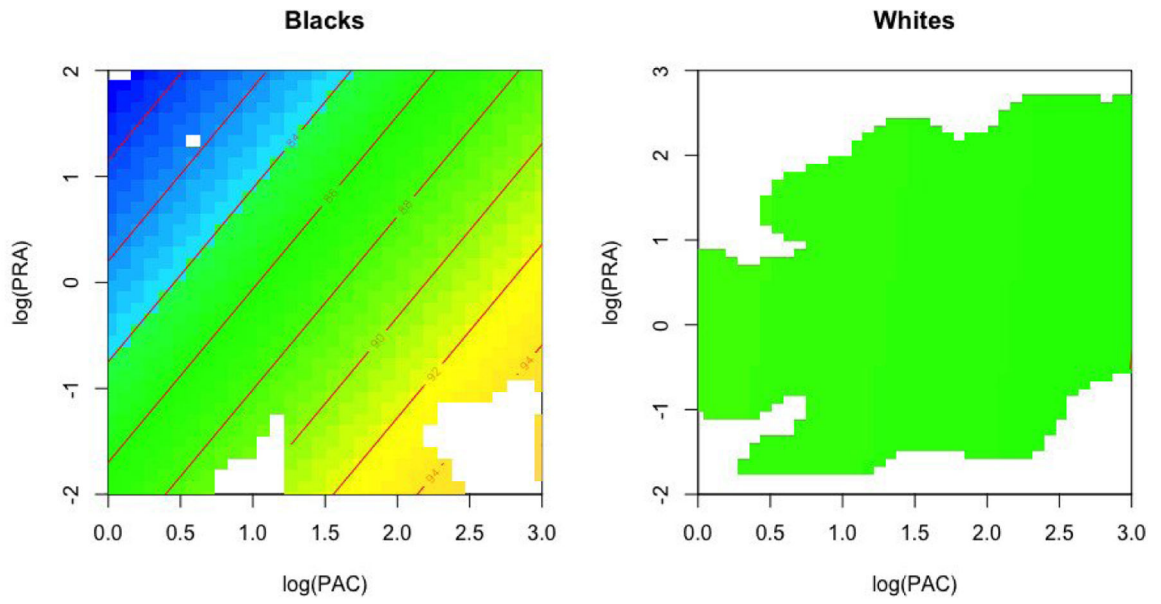
- [33]. McMurry TL, Politis DN, Banded and tapered estimates for autocovariance matrices and the linear process bootstrap, *Journal of Time Series Analysis* 31 (6) (2010) 471–482.
- [34]. Tu W, Eckert GJ, DiMeglio LA, Yu Z, Jung J, Pratt JH, Intensified effect of adiposity on blood pressure in overweight and obese children, *Hypertension* 58 (5) (2011) 818–824. [PubMed: 21968752]
- [35]. Tu W, Eckert GJ, Hannon TS, Liu H, Pratt LM, Wagner MA, DiMeglio LA, Jung J, Pratt JH, Racial differences in sensitivity of blood pressure to aldosterone, *Hypertension* 63 (6) (2014) 1212–1218. [PubMed: 24711519]
- [36]. Tu W, Eckert GJ, Pratt JH, Jan Danser A, Plasma levels of prorenin and renin in blacks and whites: their relative abundance and associations with plasma aldosterone concentration, *American Journal of Hypertension* 25 (9) (2012) 1030–1034. [PubMed: 22695510]
- [37]. Tu W, Eckert GJ, Decker BS, Howard Pratt J, Varying influences of aldosterone on the plasma potassium concentration in blacks and whites, *American Journal of Hypertension* 30 (5) (2017) 490–494. [PubMed: 28338830]



**Figure 1:**  
(a) Weight growth by race and sex; (b) Estimated weight growth curves with pointwise 95% CI, by race and sex



**Figure 2:**  
Estimated contour plots of weight as a function of height and age, by race and sex



**Figure 3:**  
 Estimated contour plots of systolic blood pressure as a function of logarithmic transformed plasma renin activity (PRA) and plasma aldosterone concentration (PAC), by race

Summary of the existing methods

**Table 1:**

Author(s), Methods	Same x(s)	Correlation	>2 Groups	Curve/Surface	Additional Comments
Bowman (2006)[17]:	Y	N	N	Curve/Surface	(+) Simple to implement and understand as a derivation from ANOVA test; (-) Assume equal variance across groups.
Dette & Neumeyer (2001)[19], Pardo-Fernandez & Van Keilegom (2007)[21], Wang & Ye (2010)[23]:	N	N	Y	Curve/Surface	(+) Demonstrated asymptotic normality of all three kernel-based statistics under $H_0$ ; Recommended wild bootstrap when studying finite samples.
Zhang & Lin (1998)[24]:	Y	Y	N	Curve	(+) Spline-based semiparametric additive model; (-) $\chi^2$ approximation can be biased with different covariate values.
Wang & Ye (2010)[23]:	N	Y	Y	Curve/Surface	(+) Able to adjust for spatial correlation; (-) Larger bias in estimating regression surface hence decreased power.
Kulasekera (1995)[13]:	N	N	N	Curve	(+) Low computational demand; (-) Low power when curve functions are similar.
Park, Hannig, & Kang (2014)[18]:	N	N	Y	Curve	(+) A visualization tool to present differences between curves across multiple locations and scales; (-) Type I error rate below nominal level; relatively low power

**Table 2:**

Average run times and standard errors in seconds over 20 runs of the proposed methods for each of the examples in the package, with and without using the parallel computing

<b>Data</b>	<b>Functions</b>	<b>Obs per group</b>	<b>Time (parallel)</b>	<b>Time (no parallel)</b>
cross-sectional	curve	(474,465)	9.6(0.2)	9.1(0.2)
	surface	(474,465)	14.8 (0.2)	23.5(0.6)
correlated	curve	(1873,1713)	86.4(0.6)	204.5(1.6)
	surface	(1873,1713)	664.6(13.6)	1749.0 (30.3)



**Table 3:**

P-values for weight growth curves and weight growth surfaces for different race and sex groups

Endpoints vs predictor(s)	Group effect	Subset data	$T_{spline}$	$T_4$	$T_3$	$T_2$
Weight vs. Age	Sex	Black	0.08	0.16	0.55	<0.01
		White	0.48	0.18	0.69	0.28
	Race	Boys	0.41	0.09	0.69	0.2
		Girls	<0.01	0.01	0.96	<0.01
Weight vs. Age & Height	Sex	Black	0.16	0.65	0.66	0.49
		White	0.49	0.55	0.55	0.77
	Race	Boys	0.42	0.42	0.46	0.35
		Girls	0.34	0.19	0.91	0.06

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript