REGULAR ARTICLE

Simplicial band depth for multivariate functional data

Sara López-Pintado · Ying Sun · Juan K. Lin · Marc G. Genton

Received: 26 April 2013 / Revised: 26 January 2014 / Accepted: 10 February 2014 / Published online: 5 March 2014 © Springer-Verlag Berlin Heidelberg 2014

Abstract We propose notions of simplicial band depth for multivariate functional data that extend the univariate functional band depth. The proposed simplicial band depths provide simple and natural criteria to measure the centrality of a trajectory within a sample of curves. Based on these depths, a sample of multivariate curves can be ordered from the center outward and order statistics can be defined. Properties of the proposed depths, such as invariance and consistency, can be established. A simulation study shows the robustness of this new definition of depth and the advantages of using a multivariate depth versus the marginal depths for detecting outliers. Real data examples from growth curves and signature data are used to illustrate the performance and usefulness of the proposed depths.

Keywords Band depth · Functional boxplot · Functional and image data · Modified band depth · Multivariate · Simplicial

Mathematics Subject Classification 62F07 · 62M10

S. López-Pintado Department of Biostatistics, Columbia University, New York, NY 10032, USA e-mail: sl2929@columbia.edu

Y. Sun Department of Statistics, The Ohio State University, Columbus, OH 43210, USA e-mail: sunwards@stat.osu.edu

J. K. Lin SearchForce, Inc., San Mateo, CA 94403, USA e-mail: juan.k.lin@gmail.com

M. G. Genton (⊠) CEMSE Division, King Abdullah University of Science and Technology, Thuwal 23955-6900, Saudi Arabia e-mail: marc.genton@kaust.edu.sa

1 Introduction

The complexity and abundance of data in emerging research fields require the improvement of statistical methodologies for analyzing complex data. In functional data analysis, each observation is a function, $x_i(t)$, i = 1, ..., n, $t \in \mathcal{I}$, where \mathcal{I} is an interval in \mathbb{R} . There are several motivations for studying functional data (see, e.g., Ramsay and Silverman 2005; Ferraty and Vieu 2006). In many research areas (e.g., medicine, biology, economics, or engineering), the data-generating process is naturally a stochastic function. Although the generating process can be a continuous function, the data are observed discretely in practice. Considering the data as functions most accurately represents the true structure of the data. In addition, if the grid on which curves are observed differs across subjects, a multivariate approach that implicitly assumes a common grid would not be valid, and it is therefore necessary to smooth the data and treat them as continuous functions defined on a common interval. Since the dimension of the observations is often significantly higher than the number of curves observed, even if the data are observed at the same time points, a standard multivariate analysis might not be computationally feasible due to the curse of dimensionality.

Often, the data recorded in studies are variables measured over time that can be considered as functional data. Classical statistical methods such as principal components analysis and regression models have recently started to be adapted to functional data. However, there are still numerous fundamental problems remaining to be investigated. In many applications, several correlated functions, say p, are observed for each sample subject. Examples include multiple lead recordings from the electrocardiogram of a patient, and weight and height paths over time for each individual. In this setting where the data are functions taking values in a multivariate space, $\mathbf{X} : \mathcal{I} \longrightarrow \mathbb{R}^p$, outliers are very difficult to detect but can affect the statistical results in many different ways. Furthermore, a multivariate outlier clearly need not be an outlier from a marginal point of view, and vice versa. The statistical analysis of multivariate curves can be significantly improved using robust estimators.

In this paper, we develop new methods for ordering multivariate functional data and for detecting outliers. A natural tool to analyze these functional data aspects is the idea of statistical depth. The notion of depth has already been used for ordering univariate functional data and it provides a measure of the "centrality" or the "outlyingness" of an observation with respect to a given data set or a population distribution (e.g., Fraiman and Muniz 2001; López-Pintado and Romo 2007, 2009; López-Pintado and Jörnsten 2007; Cuevas et al. 2007; Sun and Genton 2011, 2012a,b; Gervini 2012). We propose here an extension of the (univariate) functional band depth in López-Pintado and Romo (2009) to multivariate functional data (or trajectories). Based on this new depth, a sample of multivariate curves (or trajectories) can be ordered from the center outward and order statistics can be defined. This ordering will be a building block for extending robust statistical methods to multivariate functional data.

This paper has emerged while analyzing the early-life human growth curves using data from the 1988 National Maternal and Infant Health Survey (NMIHS) and its 1991 Longitudinal Follow-up. These data include heights and weights of boys over time as represented in Fig. 1. A preprocessing smoothing step described in López-Pintado and Wei (2011) was used to infer these smooth growth curves from discrete



Fig. 1 A three-dimensional plot of human growth curves. A representative curve is plotted in *dashed black line*

observations. We propose ways for ordering these multivariate curves (height and weight paths) and for detecting outliers. An individual might not have an unusual height or weight curve if we study these curves independently, but if we consider the height and weight curves together as a multivariate function they could show an unusual behavior.

The paper is organized as follows. In Sect. 2, the new simplicial band depth and its modified version are defined. The properties of these depths, such as consistency, are also introduced and discussed. Section 3 shows the performance of the multivariate functional depth based on simulated multivariate curves from different contaminated models. In Sect. 4, real data examples are studied. The paper ends with a discussion in Sect. 5. Proofs of theoretical results are relegated to the Appendix.

2 Simplicial band depth

2.1 Definitions

Let $\mathbf{X} : \mathcal{I} \longrightarrow \mathbb{R}^p$ be a stochastic function taking values in the space $C(\mathcal{I}, \mathbb{R}^p)$ of real continuous functions defined from a compact interval \mathcal{I} to \mathbb{R}^p and with probability distribution $P_{\mathbf{X}}$. The simplicial band depth (*SBD*) for a given function \mathbf{x} is defined as

$$SBD(\mathbf{x}, P_{\mathbf{X}}) = P\{\mathbf{x}(t) \in simplex\{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}, \forall t \in \mathcal{I}\},$$
(1)

where simplex { $X_1(t), ..., X_{p+1}(t)$ } is a random simplex in \mathbb{R}^p defined by $X_1(t), ..., X_{p+1}(t)$, which are independent copies of X evaluated at *t*. Intuitively, *SBD* measures



Fig. 2 A triangular tube determined by three black bivariate curves is shown. The curve represented in *red* (r) is completely contained in this triangular tube, while the curve represented in *blue* (b) is completely outside the tube

the probability that the trajectory of the function **x** is inside a random region in \mathbb{R}^{p+1} determined by random simplices at each time *t*.

Let $I{A}$ denote the indicator function that takes the value 1 if the event A is satisfied and zero otherwise. The sample simplicial band depth, SBD_n , of a given function **x** with respect to a sample $\mathbf{x}_1, \ldots, \mathbf{x}_n$ is:

$$SBD_n(\mathbf{x}) = \frac{1}{\binom{n}{p+1}} \sum_{1 \le i_1 < \dots < i_{p+1} \le n} I\{\mathbf{x}(t) \in \operatorname{simplex}\{\mathbf{x}_{i_1}(t), \dots, \mathbf{x}_{i_{p+1}}(t)\}, \quad \forall t \in \mathcal{I}\}.$$

 SBD_n therefore counts the proportion of regions determined by the connected simplices (convex hulls) over time that contain **x**. This provides a criterion to order the sample of multivariate curves from the center outward. In particular, for p = 2, the simplicial region is a three-dimensional (3D) tube formed by connected triangles. In this case, $SBD_n(\mathbf{x})$ measures the proportion of triangular tubes determined by three curves from the sample that contain the curve **x**. In Fig. 2, we show the triangular region determined by three bivariate curves based on a toy example. This region is constructed by connecting the simplices (triangles) determined by three given curves over time. The curve represented in red (r) is completely contained in this triangular tube and the curve in blue (b) is completely outside.

This notion of depth is quite strict, and it is often difficult for a curve to be completely contained in a triangular tube. We relax the strict containment requirement and define a modified version of the SBD at x as

$$MSBD(\mathbf{x}, P_{\mathbf{X}}) = E(\lambda[t \in \mathcal{I}, \text{ s.t. } \mathbf{x}(t) \in \text{simplex}\{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}]), \quad (2)$$

where λ is the Lebesgue measure on \mathcal{I} divided by the length of the interval \mathcal{I} . For simplicity, we assume that the length of the interval is one. Intuitively, this depth measures for how long the trajectory of $\mathbf{x}(t)$ is contained in the simplicial region determined by $\mathbf{X}_1(t), \ldots, \mathbf{X}_{p+1}(t)$. More concretely, it measures on average, the proportion of time *t* that the curve $\mathbf{x}(t)$ is in the simplicial region. If $\mathbf{x}(t)$ is completely inside the simplicial band, then $SBD(\mathbf{x}, P_{\mathbf{X}}) = MSBD(\mathbf{x}, P_{\mathbf{X}}) = 1$.

Recall that given a multivariate vector \mathbf{y} in \mathbb{R}^p and a multivariate distribution $P_{\mathbf{Y}}$, the standard simplicial depth of \mathbf{y} with respect to the distribution $P_{\mathbf{Y}}$, $SD(\mathbf{y}; P_{\mathbf{Y}})$, is defined as

$$SD(\mathbf{y}; P_{\mathbf{Y}}) = P\{\mathbf{y} \in \operatorname{simplex}\{\mathbf{Y}_1, \dots, \mathbf{Y}_{p+1}\}\},\tag{3}$$

where $\mathbf{Y}_1, \dots, \mathbf{Y}_{p+1}$ are p + 1 independent copies of \mathbf{Y} (see Liu 1990). By Fubini's theorem one can interchange the integrals in *MSBD* and express

$$MSBD(\mathbf{x}, P_{\mathbf{X}}) = \int_{\mathcal{I}} SD(\mathbf{x}(t); P_{\mathbf{X}(t)}) dt, \qquad (4)$$

where $SD(\mathbf{x}(t); P_{\mathbf{X}(t)})$ is the standard multivariate simplicial depth of \mathbf{x} at time t. Note that alternative notions of depth for multivariate functional data could be defined by just using other multivariate depths at each time t. This could be considered for future work and it is mentioned in the discussion of the paper.

The sample modified simplicial band depth, $MSBD_n(\mathbf{x})$, with respect to a sample $\mathbf{x}_1, \ldots, \mathbf{x}_n$ is:

$$MSBD_n(\mathbf{x}) = \frac{1}{\binom{n}{p+1}} \sum_{1 \le i_1 < \dots < i_{p+1} \le n} \lambda\{t \in \mathcal{I}, \text{ s.t. } \mathbf{x}(t) \in \operatorname{simplex}\{\mathbf{x}_{i_1}(t), \dots, \mathbf{x}_{i_{p+1}}(t)\}\},$$

or equivalently, using the alternative definition in (4),

$$MSBD_{n}(\mathbf{x}) = \int_{\mathcal{I}} SD_{n}(\mathbf{x}(t))dt,$$
(5)

where $SD_n(\mathbf{y})$ is the sample standard multivariate simplicial depth, with respect to a sample $\mathbf{y}_1, \ldots, \mathbf{y}_n$ and can be expressed as

$$SD_n(\mathbf{y}) = \frac{1}{\binom{n}{p+1}} \sum_{1 \le i_1 < \dots < i_{p+1} \le n} I\{\mathbf{y} \in \operatorname{simplex}\{\mathbf{y}_{i_1}, \dots, \mathbf{y}_{i_{p+1}}\}\}.$$
 (6)

Note that for the case of p = 2, $MSBD_n$ measures, on average, the proportion of times *t* where the curve **x** is inside a triangular tube determined by three curves from the sample. If p = 1, then definitions in Eqs. (1) and (2) reduce to the notions of band depth and modified band depth introduced by López-Pintado and Romo (2009) for univariate functional data. Computations of SBD_n and $MSBD_n$ can be based on the fast algorithm for computing band depth of Sun et al. (2012) combined with fast algorithms to compute simplicial depth (see, e.g., Rousseeuw and Ruts 1996; Cheng and Ouyang 2001).

2.2 Properties

In this section, we describe some properties satisfied by *SBD* and *MSBD*. They can be derived from the band depth, modified band depth and simplicial depth properties in López-Pintado and Romo (2009) and Liu (1990), respectively. The proofs of the properties are given in the Appendix.

We state some assumptions needed for the properties. Let $\mathbf{X} : \mathcal{I} \longrightarrow \mathbb{R}^p$ be a stochastic function taking values in the space $C(\mathcal{I}, \mathbb{R}^p)$ of continuous functions, and with probability distribution $P_{\mathbf{X}}$. Also assume that \mathcal{I} is a compact interval, that the probability $P_{\mathbf{X}}$ is absolutely continuous and $P_{\mathbf{X}(t)}$ has a unique deepest point at each t. In other words, there is one unique point that maximizes the multivariate simplicial depth, $SD(\mathbf{x}(t); P_{\mathbf{X}(t)})$, at each time t.

Theorem 1 Under the previous assumptions we can state the following properties for SBD:

- 1. The simplicial band depth is invariant under the following transformations:
- (a) Let $\mathbf{T}(\mathbf{x})$ be the combined function defined as $\mathbf{T}(\mathbf{x}(t)) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)$, where $t \in \mathcal{I}$ and $\mathbf{A}(t)$ is a $p \times p$ invertible matrix with $A_{ij}(t)$ a continuous function in t with $t \in \mathcal{I}$ and $\mathbf{b}(t) \in C(\mathcal{I}, \mathbb{R}^p)$. For simplicity, we call these assumptions "standard assumptions for the linear transformation" for later use in the paper. Then,

$$SBD(\mathbf{T}(\mathbf{x}), P_{T(X)}) = SBD(\mathbf{x}, P_X).$$

(b) Let g be a one-to-one transformation of the interval \mathcal{I} . Then,

$$SBD(\mathbf{x}(g), P_{\mathbf{X}(g)}) = SBD(\mathbf{x}, P_{\mathbf{X}}).$$

2. Vanishing at infinity: $SBD(\mathbf{x}, P_{\mathbf{X}})$ converges to zero when the supremum norm of the components of the multivariate process \mathbf{x} tends to infinity:

 $\sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} SBD(\mathbf{x}, P_{\mathbf{X}}) \longrightarrow 0, \quad when \ M \to \infty,$

where $||x_k||_{\infty}$ is the supremum norm of the kth component of **x**.

Theorem 2 Under the assumptions stated at the beginning of this section, MSBD satisfies the following properties:

1. Invariance of MSBD: Let $\mathbf{T}(\mathbf{x})$ be the combined function defined as $\mathbf{T}(\mathbf{x}(t)) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)$. Under the standard assumptions for the linear transformation in Theorem 1, we can show that:

$$MSBD(\mathbf{T}(\mathbf{x}), P_{T(\mathbf{X})}) = MSBD(\mathbf{x}, P_{\mathbf{X}}).$$

2. Monotonicity with respect to the deepest point: For any $c \in [0, 1]$ we have that

$$MSBD(\mathbf{x}, P_{\mathbf{X}}) \leq MSBD(\mathbf{y} + c(\mathbf{x} - \mathbf{y}), P_{\mathbf{X}}),$$

where **y** is the deepest point based on $MSBD(\mathbf{x}, P_{\mathbf{X}})$.

3. Maximality at the center:

$$MSBD(\mathbf{z}, P_{\mathbf{X}}) = \sup_{\mathbf{x} \in C(\mathcal{I}, \mathbb{R}^p)} MSBD(\mathbf{x}, P_{\mathbf{X}}),$$

for any distribution $P_{\mathbf{X}}$ with unique center function of symmetry \mathbf{z} .

4. Vanishing at infinity: $MSBD(\mathbf{x}, P_{\mathbf{X}})$ converges to zero when the norm of the multivariate process at each component x_k , with $1 \le k \le p$ tends to infinity for almost all time points in \mathcal{I} . Specifically, for any sequence of functions \mathbf{x}_n from $C(\mathcal{I}, \mathbb{R}^p)$ then $\lim_{n\to\infty} MSBD(\mathbf{x}_n, P_{\mathbf{X}}) = 0$ if $\lim_{n\to\infty} |x_{n,k}(t)| = \infty$ for almost all time points t in \mathcal{I} , where $1 \le k \le p$ and $x_{n,k}$ is the kth component of the multivariate function \mathbf{x}_n .

2.3 Finite-dimensional case

In practice, the full curves are not observed. Instead curves are evaluated at a finite set of time points, t_1, t_2, \ldots, t_k in \mathcal{I} . Under this assumption one can still define *MSBD* as in (2) considering a finite measure that counts the proportion of time points satisfying the corresponding condition. In this setting, $\mathbf{X} = {\mathbf{X}(t_1), \mathbf{X}(t_2), \ldots, \mathbf{X}(t_k)}$ is a random vector in $(\mathbb{R}^p)^k$ and at each given time point $\mathbf{X}(t)$ is a *p*-dimensional random vector with a given distribution function $P_{\mathbf{X}(t)}$. In particular, the finite dimensional definitions of *MSBD* and *MSBD*_n are

$$MSBD(\mathbf{x}, P_{\mathbf{X}}) = E\left(\frac{1}{k}\sum_{1\leq j\leq k} I\left[\mathbf{x}(t_{j}) \in \operatorname{simplex}\{\mathbf{X}_{1}(t_{j}), \dots, \mathbf{X}_{p+1}(t_{j})\}\right]\right)$$
$$= \frac{1}{k}\sum_{1\leq j\leq k} P(\mathbf{x}(t_{j}) \in \operatorname{simplex}\{\mathbf{X}_{1}(t_{j}), \dots, \mathbf{X}_{p+1}(t_{j})\}),$$
$$MSBD_{n}(\mathbf{x}) = \frac{1}{\binom{n}{p+1}}\sum_{1\leq i_{1}<\dots< i_{p+1}\leq n, 1\leq j\leq k}\sum_{i_{1}\leq i_{1}\leq i_{1}\leq i_{2}\leq k}\frac{1}{k}I\left[\mathbf{x}(t_{j})\right]$$
$$\in \operatorname{simplex}\{\mathbf{x}_{i_{1}}(t_{j}), \dots, \mathbf{x}_{i_{p+1}}(t_{j})\}\right].$$

Theorem 3 (Consistency of $MSBD_n$): In this finite-dimensional setting and under the assumptions stated at the beginning of this section, the sample $MSBD_n$ converges uniformly almost surely to the population MSBD as $n \to \infty$:

$$\sup_{\mathbf{x}\in(\mathbb{R}^p)^k}|MSBD_n(\mathbf{x})-MSBD(\mathbf{x}, P_{\mathbf{X}})|\xrightarrow{a.s.} 0.$$

3 Monte Carlo simulation study

Since the order of the multivariate functional data is determined by depth measure, different data depths usually lead to different robust estimators of the distribution

function that can be used to construct outlier detection rules (Sun and Genton 2011, 2012a). Hence, the assessment of outlier detection provides a way to compare different notions of data depth. In this section, we perform a simulation study comparing the performance of the *MSBD* to the marginal modified band depth (*MBD*), and to an alternative notion of multivariate depth, called weighted modified band depth (*WMBD*), that has been recently proposed by Ieva and Paganoni (2013). The *WMBD* consists of taking a weighted average of the marginal modified band depths of each component of the multivariate function. We consider three different settings of the weights that are required to be pre-specified in Ieva and Paganoni (2013). The reason for choosing *MSBD* instead of *SBD* is that the depth values induced by *SBD* are very likely to have ties in practice, and effective tie breaking techniques are needed in order to produce a meaningful ranking.

For data generation, we consider bivariate curves, $\mathbf{X}_i(t) = (X_{1i}(t), X_{2i}(t))^T$, i = 1, ..., n, generated from different models introducing outliers to $X_1(t)$ and $X_2(t)$. The simulation designs are similar to those in López-Pintado and Romo (2009) and Sun and Genton (2011), but aimed at generating bivariate functional data. Model 1 is a basic one without contamination. Models 2, 3 and 4 have magnitude outliers, i.e., curves that are shifted away from the bulk of the data. Model details are described as follows:

1. Model 1 is $\mathbf{X}_i(t) = \mathbf{Z}_i(t), t \in [0, 1]$, where $\mathbf{Z}_i(t)$ is a stochastic bivariate Gaussian process with zero mean and a bivariate Matérn cross-covariance function (Gneiting et al. 2010; Apanasovich et al. 2012):

$$C_{ii}(s,t) = \sigma_{ii}M(|s-t|;v_{ii},\alpha_{ii}), \quad i, j = 1, 2,$$

where $M(h; \nu, \alpha) = \frac{\sigma^2}{2^{\nu-1}\Gamma(\nu)} (\alpha h)^{\nu} \mathcal{K}_{\nu}(\alpha h), h = |s - t| \in [0, 1]$, is a Matérn class (Matérn 1960; Stein 1999), \mathcal{K}_{ν} is a modified Bessel function of the second kind, σ^2 is the marginal variance, $\nu > 0$ is a smoothness parameter, and $\alpha > 0$ is a scale parameter;

- 2. Model 2 includes contamination for $X_1(t)$ and $X_2(t)$: $\mathbf{X}_i(t) = c_i K \mathbf{Z}_i(t)$, where c_i is 1 with probability q and 1/K with probability 1 q, K is a contamination size constant;
- 3. Model 3 is partially contaminated: $\mathbf{X}_i(t) = c_i K \mathbf{Z}_{1i}(t)$, if $t \ge T_i$ and $\mathbf{X}_i(t) = \mathbf{Z}_i(t)$, if $t < T_i$, where T_i is a random number generated from a uniform distribution on [0, 1];
- 4. Model 4 is contaminated by peaks: $\mathbf{X}_i(t) = c_i K \mathbf{Z}_i(t)$, if $T_i \leq t \leq T_i + \ell$, and $\mathbf{X}_i(t) = \mathbf{Z}_i(t)$ otherwise, where T_i is a random number from a uniform distribution in $[0, 1 \ell]$.

Following Sun and Genton (2011, 2012a), we use the empirical rule of a constant factor *F* times the 50% central region in the adjusted functional boxplot to detect the outliers in $X_1(t)$, where the 50% central region contains the 50% observations with the largest depth values. We first generate n = 100 curves evaluated at 50

equally spaced time points with parameters q = 0.1, K = 12, $\ell = 10/49$. The parameters in the bivariate Matérn model are $\sigma_{11} = \sigma_{22} = 1$, $\alpha_{11} = 0.02$, $\alpha_{22} = 0.01$, $\nu_{11} = 1.2$, $\nu_{22} = 0.6$, $\alpha_{12} = 0.016$, $\nu_{12} = 1.0$. At a given time t, $\mathbf{X}(t)$ follows a bivariate normal distribution with correlation $\rho = 0.1$, 0.3, 0.6. Then we compute depth values using bivariate *MSBD*, marginal *MBD*, and *WMBD* proposed in Ieva and Paganoni (2013), where for the latter we consider three different sets of weights, (0.25, 0.75), (0.5, 0.5) and (0.75, 0.25), and denote these depths as *WMBD1*, *WMBD2* and *WMBD3*. For the *MBD* and *WMBD*, the outliers in $X_1(t)$ are detected by the adjusted functional boxplot, where the 50% central region of $X_1(t)$ is defined using the ordering induced by the corresponding depth. For the bivariate *MSBD*, the 50% central region is in 3D and the constant factor F is defined as the orthogonal distance to the 45° line.

As in Sun and Genton (2011), we assess outlier detection performance by examining the distribution of two quantities: \hat{p}_c , the percentage of correctly detected outliers (the number of correctly detected outliers divided by the total number of outlying curves), and \hat{p}_f , the percentage of falsely detected outliers (the number of falsely detected outliers divided by the total number of non-outlying curves). We run 1,000 replications. For model 1, with no outliers, we estimate the percentage, \hat{p}_0 , of times that the different methods detect no outliers, and also the percentage \hat{p}_f . For models 2, 3 and 4, we estimate the percentages \hat{p}_c and \hat{p}_f . Table 1 shows the means and standard deviations of these quantities. Good outlier detection performance is given by high correct detection percentages, \hat{p}_0 and \hat{p}_c , and a low false detection percentage, \hat{p}_f . The factor F (shown in Table 1) in the adjusted functional plot is chosen to make $\hat{p}_0 = 99.3\%$ (Sun and Genton 2012a), for each combination of the five depth notions and the three values of ρ . For the marginal *MBD*, F remains the same while it varies with ρ for the *MSBD* and *WMBD*. The selected F also depends on the weights of *WMBD*. Therefore, we compare the outlier detection performance for a given ρ . For all these cases, the bivariate *MSBD* clearly outperforms the marginal MBD and the WMBD. In particular, for all methods but WMBD3 the percentage of falsely detected outliers (\hat{p}_f) is close to zero. The main difference in performance is in the percentage of correctly detected outliers (\hat{p}_c) which is clearly higher using MSBD than any other method for all models. Intuitively, if the contamination leads to a different direction of the correlation while the corresponding marginal contamination is not large enough to be detected (see Sun and Genton 2011, 2012a), the bivariate MSBD should perform better.

For the *WMBD*, the outlier detection performance depends on the weights which need to be pre-specified. However, there is no justification for optimal weights making this method impractical. In these simulations, the performance of *WMBD* is similar to the marginal *MBD* and both are worse than the bivariate *MSBD*. Similar to the usual bivariate case, bivariate functional outliers in bivariate functional data are not necessarily marginal outliers. The outlier detection method based on the marginal *MBD* ordering therefore fails to detect such type of outliers, and a multivariate functional depth is recommended.

Table 1 The mean and standard deviation (in parentheses) of the percentages \hat{p}_c and \hat{p}_f for marginal *MBD* and bivariate *MSBD*, and for *WMBD* with different weights, with $\rho = 0.1, 0.3, 0.6, 1,000$ replications, n = 100 curves for models 1–4, and adjustment factor *F*

ρ	Model 1		Model 2		Model 3		Model 4	
	F	\hat{p}_f	\hat{p}_c	\hat{p}_f	$\overline{\hat{p}_c}$	\hat{p}_f	\hat{p}_c	\hat{p}_f
0.1								
MSBD	1.68	0.01 (0.08)	91.5 (9.4)	0.00 (0.06)	90.5 (12.3)	0.00 (0.07)	88.6 (15.8)	0.00 (0.07)
MBD	3.00	0.01 (0.03)	67.9 (15.9)	0.00 (0.03)	66.9 (15.6)	0.00 (0.00)	66.8 (15.6)	0.00 (0.00)
WMBD1	0.75	0.01 (0.10)	66.1 (18.3)	0.00 (0.05)	62.4 (21.5)	0.00 (0.03)	60.7 (22.5)	0.00 (0.03)
WMBD2	1.46	0.01 (0.10)	70.1 (15.6)	0.00 (0.06)	69.0 (15.5)	0.00 (0.05)	68.8 (15.8)	0.00 (0.05)
WMBD3	2.30	0.01 (0.10)	79.0 (14.0)	0.19 (0.51)	69.6 (15.0)	0.00 (0.06)	69.6 (15.1)	0.00 (0.06)
0.3								
MSBD	1.97	0.01 (0.08)	90.8 (10.0)	0.00 (0.06)	89.6 (12.8)	0.00 (0.07)	87.6 (16.6)	0.01 (0.08)
MBD	3.00	0.01 (0.03)	67.9 (15.9)	0.00 (0.03)	66.9 (15.6)	0.00 (0.00)	66.8 (15.6)	0.00 (0.00)
WMBD1	0.74	0.01 (0.08)	67.2 (18.7)	0.01 (0.12)	64.2 (20.8)	0.00 (0.05)	63.0 (21.5)	0.00 (0.05)
WMBD2	1.50	0.01 (0.08)	70.3 (15.5)	0.01 (0.08)	69.0 (15.7)	0.00 (0.00)	68.9 (15.8)	0.00 (0.00)
WMBD3	2.37	0.01 (0.08)	78.7 (14.1)	0.16 (0.44)	69.0 (15.2)	0.00 (0.05)	69.0 (15.1)	0.00 (0.06)
0.6								
MSBD	3.00	0.01 (0.08)	87.0 (11.4)	0.00 (0.07)	85.4 (14.8)	0.00 (0.06)	83.3 (18.7)	0.00 (0.06)
MBD	3.00	0.00 (0.03)	67.9 (15.9)	0.00 (0.03)	66.9 (15.6)	0.00 (0.00)	66.8 (15.6)	0.00 (0.00)
WMBD1	0.90	0.01 (0.13)	68.4 (18.7)	0.01 (0.12)	66.5 (19.4)	0.02 (0.16)	64.1 (21.4)	0.02 (0.17)
WMBD2	1.70	0.01 (0.08)	69.2 (15.5)	0.00 (0.04)	68.3 (15.5)	0.00 (0.03)	68.2 (15.7)	0.00 (0.03)
WMBD3	2.45	0.01 (0.10)	77.1 (14.4)	0.08 (0.32)	68.5 (15.4)	0.00 (0.05)	68.5 (15.3)	0.00 (0.05)

4 Data examples

4.1 Chinese script data

The first data set consists of repeated writings of Chinese words by the same person (see Fig. 3). This data set was obtained from Ramsay et al. (2009) and had already been preprocessed and registered. The Chinese script data can be seen as bivariate functional data, where each script is a trajectory over time. In Fig. 3, we represent the projection of the data in the X-Y axis, although these values are parametrized by time, which is the third dimension. The multivariate simplicial band depth provides a ranking of these scripts from the center outward. Based on this ranking, the median or most representative script within the sample can be defined and possible outliers can be detected. In Fig. 3, we show the deepest trajectory (in black) and the triangular tube determined by the three deepest curves. It can be seen that the deepest signature is representative of the sample of scripts and can be used as a median trajectory. The ordering provided by the simplicial band depth can be used as a building block for generalizing different robust statistical methods to multivariate functional data.



Fig. 3 Chinese script replicated 100 times. The triangular tube determined by the three deepest *curves* is shown. *Black dots* indicate observations along the deepest curve

4.2 Growth curves data

The data in this second example, as mentioned in the introduction, consist of the height and weight over time for a sample of infants from birth to two years of age. The data are sparse so a preprocessing step is needed to estimate the curves on a common set of time points (see López-Pintado and Wei 2011). The study is motivated by an early-life human growth project using data from the 1988 National Maternal and Infant Health Survey (NMIHS) and its 1991 Longitudinal Follow-up. The study included 2,555 boys and 2,510 girls born in the US in the calendar year of 1988. Their heights and weights were taken sporadically only when they visited a hospital. Consequently, their growth paths were observed on small, variable sets of sparse and irregularly spaced time points. The study of growth patterns of infants has long been an important research topic in epidemiology. The most informative growth pattern is represented by the underlying height and weight processes as continuous functions of age. The simplicial band depth provides a way of ordering these bivariate curves from the sample, and a median growth curve can be defined. In addition, the notion of depth can be used to detect outliers, which in this case correspond to unusual growth patterns. In Fig. 4, the raw height and weight data are represented for a subset of 150 boys. Before computing the depths, the data were smoothed using the approach in López-Pintado and Wei (2011).

In Fig. 5 these smoothed curves are represented in separate 2-dimensional plots. We have compared the ordering provided by the marginal depths with the ordering given by the multivariate simplicial band depth. The ordering of the height curves, X(t), changes when using the multivariate simplicial band depth instead of the univariate marginal depth since we are incorporating information about the boys' weights, Y(t).



Fig. 4 Sparse growth data. The height (*left panel*) and weight (*right panel*) over time for a sample of boys are represented

In Fig. 5a, b, the deepest bivariate curve using the simplicial band depth, *SBD*, is represented in solid (red) line and the deepest curves using the marginal univariate band depths for height and weight are represented with a dashed (black) curve. The marginal deepest curves correspond to two different boys although their growth curve is close to the deepest one obtained using the multivariate depth. This implies that the deepest bivariate height-weight curve shows similar behaviour to the deepest marginal curves.

Another main application of simplicial band depth is the detection of multivariate outliers that are not necessarily marginal outliers. To illustrate this, in the second row of Fig. 5 we represent in a solid (red) line the smoothed height and weight curves for a boy from the sample. This bivariate curve is the second percentile based on the center-outward order provided by the multivariate simplicial band depth, whereas it is the 25th and 29th percentile if we apply the marginal band depth to weight and height, respectively. It can be seen that for this individual, the height growth curve increases very steadily (especially after one year of age) whereas the weight curve increases rapidly over time. Looking at the height and weight curves, we can see that the child is initially at a very high percentile in height and a low percentile in weight. As the child grows, the percentile in weight starts increasing while the percentile in height decreases. This could be a sign of a weight problem and it is important to detect this problem as early as possible. Marginally, these height and weight curves are not considered outliers, but the bivariate growth curve shows an unusual pattern when the curves are ordered according to the multivariate simplicial band depth. The multivariate simplicial band depth can be useful for screening out these types of outliers that can be indicative of a growth issue in a child.

We applied the functional boxplot ideas proposed by Sun and Genton (2011) to the growth data. The functional boxplots are useful tools for visualizing central regions and for detecting outliers. In Fig. 6, we represent the boxplots for the weight (top row) and height (bottom row) curves based on marginal *MBD* (left column) and mul-

35

30

25

(a)





Fig. 5 The height and weight *curves* are represented in the *left* and *right panels*, respectively. The deepest bivariate curve using the order provided by the multivariate simplicial band depth (SBD) is represented in a solid (red) line in (a) and (b). The deepest curves using marginal band depth (BD) for height and weight independently are represented by dashed (black) curves. In the bottom row, the thicker solid (red) curves in (c) and (d) correspond to the bivariate height-weight curve of an individual who is the second percentile using SBD ordering. This child's height curve is the 29th percentile based on marginal BD in height and the 24th percentile for BD in weight

tivariate MSBD (right column). It can be seen that the 50% central regions presented as the middle boxes in the functional boxplots are wider when using the multivariate MSBD. Since the 50% central regions contain the same number of curves, the narrower boxes indicate that some relatively deep curves of weight or height based on the marginal *MBD* are ordered as more shallow ones based on the *MSBD*, possibly due to their relatively different correlation patterns. Since the outlier detection depends on the width of the box in a functional boxplot, it is not surprising that one outlier is detected in the functional boxplot of the weight curves based on the marginal MBD while no outliers are found based on the MSBD. More investigation on this potential outlier is needed, including the height and weight measurements given the age, and the growth patterns in both height and weight. Moreover, in this example, we focus on visualization and comparisons using functional boxplots without any adjustment. However, if outlier detection is of interest, the adjusted functional boxplot (Sun and Genton 2012a) should be considered.



Fig. 6 The functional *boxplots* for the weight curves using the marginal modified band depth (*MBD*) and multivariate modified simplicial band depth (*MSBD*) are represented in (**a**) and (**b**), respectively. The functional *boxplots* for the height curves based on the marginal modified band depth (*MBD*) and the multivariate modified simplicial band depth (*MSBD*) are represented in (**c**) and (**d**), respectively

5 Discussion

The analysis of functional data is a developing field in statistics that has emerged in the last decade along with the dramatic growth of large and complex data sets. In many applications, the basic underlying observation is a multivariate function. Height and weight measurements of a child over time, and different leads in the electrocardiogram of a patient, are examples. New challenges arise when the functions are multivariate. In this paper, we have extended the notion of band and modified band depths for univariate functional data to multivariate functional data. This provides a robust ordering of multivariate functions from the center outward, and robust statistics such as the median or trimmed mean. It also provides a method for detecting outliers that are not easy to screen out when analyzing complex data. As in the standard multivariate context, a multivariate functional outlier is not necessarily an outlier in the marginal data. Multivariate orderings are therefore needed to detect multivariate outliers.

Recall that using Fubini's theorem the *MSBD* is equivalent to the integral over time of the multivariate simplicial depth; see (4). Other notions of multivariate depth, such as Tukey's depth, could have been used instead of the simplicial depth to define

alternative notions of multivariate functional depths. These multivariate depths are computationally feasible in the bivariate functional data cases that we have considered in the paper. When the multivariate functional data takes values in spaces of higher dimension other multivariate depths could be considered. For example, the finite dimensional version of the band depth proposed in López-Pintado and Romo (2009) could be used as the multivariate depth at each time since it can be computed fast in any dimension. The theoretical properties of these alternative ways of defining multivariate functional depths is an interesting topic for future research. An appealing aspect of these methods is that they can be used to construct boxplots to visualize the most representative and outlying curves in the sample using the ordering provided by the multivariate depths. All of these tools can be extended to image data. The notion of band depth has already been adapted to spatial data by Sun and Genton (2012a) and generalized to volume depth for images and surfaces by Genton et al. (2014).

We could extend the ideas proposed in our paper to multivariate images, where multiple correlated images are observed for each patient. In this case, the functions from the sample are defined from a bivariate space (pixels) to a multivariate space. In this setting, it will be important to have computationally efficient procedures to obtain the simplicial band depth.

Appendix

Proof of Theorem 1

1(a). Let $\mathbf{T}(\mathbf{x})$ be the combined function defined as $\mathbf{T}(\mathbf{x}(t)) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)$, where $t \in \mathcal{I}$. Assume it satisfies the standard assumptions for the linear transformation presented in Theorem 1. By definition,

$$SBD(\mathbf{x}, P_{\mathbf{X}}) = P\{\mathbf{x}(t) \in simplex\{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}, \forall t \in \mathcal{I}\}.$$

It is trivial to check that for any fixed $t \in \mathcal{I}$, we have that the curve $\mathbf{x}(t) \in \operatorname{simplex} \{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}$ if and only if the curve $\mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t) \in \operatorname{simplex} \{\mathbf{A}(t)\mathbf{X}_1(t) + \mathbf{b}(t), \dots, \mathbf{A}(t)\mathbf{X}_{p+1}(t) + \mathbf{b}(t)\}$ and therefore,

$$SBD(\mathbf{T}(\mathbf{x}), P_{\mathbf{T}(\mathbf{X})}) = SBD(\mathbf{x}, P_{\mathbf{X}}).$$

1(b). Let g be a one-to-one transformation of the interval \mathcal{I} . It is straightforward to prove that for any fixed $t \in \mathcal{I}$, $\mathbf{x}(g(t)) \in \operatorname{simplex}\{\mathbf{X}_1(g(t)), \ldots, \mathbf{X}_{p+1}(g(t))\}$ if and only if $\mathbf{x}(t) \in \operatorname{simplex}\{\mathbf{X}_1(t), \ldots, \mathbf{X}_{p+1}(t)\}$ and therefore,

$$SBD(\mathbf{x}(g), P_{\mathbf{X}(g)}) = SBD(\mathbf{x}, P_{\mathbf{X}})$$

 SBD(x, P_X) converges to zero when the supremum norm of the components of the process x tend to infinity. Specifically,

$$\sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} SBD(\mathbf{x}, P_{\mathbf{X}}) \longrightarrow 0, \text{ when } M \to \infty,$$

where x_k is the *k*th component of the multivariate function **x**. We establish first by contradiction the inclusion

$$\left\{ (\mathbf{X}_1, \dots, \mathbf{X}_{p+1}) \colon \mathbf{x}(t) \in \operatorname{simplex} \{ \mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t) \}, \forall t \in \mathcal{I} \right\}$$
$$\subset \bigcup_{r=1}^{p+1} \bigcup_{k=1}^p \left\{ (\mathbf{X}_1, \dots, \mathbf{X}_{p+1}) \colon \|X_{rk}\|_{\infty} \ge \|x_k\|_{\infty} \right\}$$

where X_{rk} is the *k*th component of the multivariate function \mathbf{X}_r . If $\mathbf{x}(t) \in \text{simplex}\{\mathbf{X}_1(t), \mathbf{X}_2(t), \dots, \mathbf{X}_{p+1}(t)\}$ for all $t \in \mathcal{I}$, then, for each *k* and $t \in \mathcal{I}$,

$$\min_{r=1,\dots,p+1} \{X_{rk}(t)\} \le x_k(t) \le \max_{r=1,\dots,p+1} \{X_{rk}(t)\}.$$
(7)

Assume that $||X_{r,k}||_{\infty} < ||x_k||_{\infty}$ for each k = 1, ..., p, and r = 1, ..., p + 1; this implies that, for each r and k, we have

$$\max_{t\in\mathcal{I}}|X_{rk}(t)|<\max_{t\in\mathcal{I}}|x_k(t)|.$$

Let t_k be the point where the maximum of $x_k(t)$ is achieved. Then, for all $r = 1, ..., p + 1, |X_{rk}(t_k)| < |x_k(t_k)|$, which contradicts (7). Therefore,

$$\sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} SBD(\mathbf{x}, P_{\mathbf{X}})$$

$$\leq \sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} P\left(\mathbf{x}(t) \in \operatorname{simplex}\{\mathbf{X}_1(t), \dots, \mathbf{X}_j(t)\}, \forall t \in \mathcal{I}\right)$$

$$\leq \sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} \sum_{r=1}^{p+1} \sum_{k=1}^p P\left(\|X_{rk}\|_{\infty} \ge \|x_k\|_{\infty}\right)$$

$$\leq \sum_{r=1}^{p+1} \sum_{k=1}^p \sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} P\left(\|X_{rk}\|_{\infty} \ge \|x_k\|_{\infty}\right)$$

and this implies that $\sup_{\min_{k=1,\dots,p} \|x_k\|_{\infty} \ge M} SBD(\mathbf{x}, P_{\mathbf{X}}) \longrightarrow 0$, when $M \to \infty$.

Proof of Theorem 2

1. Let $\mathbf{T}(\mathbf{x}(t)) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)$ be a linear transformation satisfying the standard assumptions in Theorem 1. By definition,

$$MSBD(\mathbf{x}, P_{\mathbf{X}}) = E(\lambda[t \in \mathcal{I}, \text{ s.t. } \mathbf{x}(t) \in \text{simplex}\{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}]).$$
(8)

It is trivial to check that for any fixed $t \in \mathcal{I}$, we have that the curve $\mathbf{x}(t) \in \operatorname{simplex}\{\mathbf{X}_1(t), \dots, \mathbf{X}_{p+1}(t)\}$ if and only if the curve $\mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t) \in \operatorname{simplex}\{\mathbf{A}(t)\mathbf{X}_1(t) + \mathbf{b}(t), \dots, \mathbf{A}(t)\mathbf{X}_{p+1}(t) + \mathbf{b}(t)\}$ and therefore,

$$MSBD(\mathbf{T}(\mathbf{x}), P_{\mathbf{T}(\mathbf{X})}) = MSBD(\mathbf{x}, P_{\mathbf{X}}).$$

- 2. The monotonicity property follows directly from expression (4) and the monotonicity property satisfied by the simplicial depth *SD* defined in Eq. (3). □
- 3. If $P_{\mathbf{X}}$ has unique center of symmetry $\mathbf{y} \in C(\mathcal{I}, \mathbb{R}^p)$, in the sense that $P_{\mathbf{X}-\mathbf{y}} = P_{\mathbf{y}-\mathbf{X}}$, then for every $t \in \mathcal{I}, \mathbf{y}(t)$ is the center of symmetry for $P_{\mathbf{X}(t)}$. Therefore, using the alternative expression of *MSBD* in (4), and since *SD* is maximized at the center, the proof is concluded.
- 4. The vanishing at infinity property follows directly from expression (4) and the vanishing at infinity properties of *SD*. □

Proof of Theorem 3 By interchanging the two sums in the definition of $MSBD_n$, one can write

$$MSBD_{n}(\mathbf{x}) = \frac{1}{k} \sum_{1 \le j \le k} \frac{1}{\binom{n}{p+1}} \sum_{1 \le i_{1} < \dots < i_{p+1} \le n,} I[\mathbf{x}(t_{j}) \in simplex\{\mathbf{x}_{i_{1}}(t_{j}), \dots, \mathbf{x}_{i_{p+1}}(t_{j})\}],$$

which is equivalent to

$$MSBD_n(\mathbf{x}) = \frac{1}{k} \sum_{1 \le j \le k} SD_n(\mathbf{x}(t_j)),$$

where $SD_n(\mathbf{x}(t_j))$ is the sample simplicial depth of $\mathbf{x}(t_j)$ as defined in Eq. (6). Also, note that one can write $MSBD(\mathbf{x}, P_{\mathbf{X}}) = \frac{1}{k} \sum_{1 \le j \le k} SD(\mathbf{x}(t_j), P_{\mathbf{X}(t_j)})$, where $SD(\mathbf{x}(t_j), P_{\mathbf{X}(t_j)})$ is the population simplicial depth of $\mathbf{x}(t_j)$ as defined in Eq. (3). Therefore,

$$\sup_{\mathbf{x}\in(\mathbb{R}^p)^k} |MSBD_n(\mathbf{x}) - MSBD(\mathbf{x}, P_{\mathbf{X}})|$$
$$= \sup_{\mathbf{x}\in(\mathbb{R}^p)^k} \left| \frac{1}{k} \sum_{1\leq j\leq k} (SD_n(\mathbf{x}(t_j)) - SD(\mathbf{x}(t_j), P_{\mathbf{X}(t_j)})) \right|$$

By the uniform consistency of the sample simplicial depth proven in Liu (1990) we can conclude that the sample $MSBD_n$ converges uniformly almost surely to the population MSBD as $n \to \infty$:

$$\sup_{\mathbf{x}\in(\mathbb{R}^p)^k}|MSBD_n(\mathbf{x})-MSBD(\mathbf{x},P_{\mathbf{X}})|\xrightarrow{a.s.}0.$$

References

Apanasovich TV, Genton MG, Sun Y (2012) A valid Matérn class of cross-covariance functions for multivariate random fields with any number of components. J Am Stat Assoc 107:180–193

- Cheng A, Ouyang M (2001) On algorithms for simplicial depth. In: Proceeding 13th Canadian conference on computational geometry, vol 1, pp 53–56
- Cuevas A, Febrero M, Fraiman R (2007) Robust estimation and classification for functional data via projection-based depth notions. Comput Stat 22:481–496
- Ferraty F, Vieu P (2006) Nonparametric Functional Data Analysis. Springer, New York

Fraiman R, Muniz G (2001) Trimmed means for functional data. Test 10:419-440

- Genton MG, Johnson C, Potter K, Stenchikov G, Sun Y (2014) Surface boxplots. Stat 3:1-11
- Gervini D (2012) Outlier detection and trimmed estimation for general functional data. Statistica Sinica 22:1639–1660
- Gneiting T, Kleiber W, Schlather M (2010) Matérn cross-covariance functions for multivariate random fields. J Am Stat Assoc 105:1167–1177
- Ieva F, Paganoni M (2013) Depth measures for multivariate functional data. Commun Stat 42(7):1265–1276 Liu RY (1990) On a notion of data depth based upon random simplices. Ann Stat 18:405–414
- López-Pintado S, Jörnsten R (2007) Functional data analysis via extensions of the band depth, IMS lecture Notes-Monograph Series. Inst Math Stat 54:103–120
- López-Pintado S, Romo J (2007) Depth-based inference for functional data. Comput Stat Data Anal 51:4957–4968
- López-Pintado S, Romo J (2009) On the concept of depth for functional data. J Am Stat Assoc 104(486): 718–734
- López-Pintado S, Wei Y (2011) Depth for sparse functional data. In: Ferraty F (ed) Recent advances in functional data analysis and related topics. Springer, Berlin, pp 209–212
- Matérn B (1960) Spatial variation. Springer, New York
- Ramsay JO, Hooker G, Graves S (2009) Functional data analysis with R and MATLAB. Springer, New York
- Ramsay JO, Silverman BW (2005) Functional data analysis, 2nd edn. Springer, New York
- Rousseeuw PJ, Ruts I (1996) Bivariate location depth. Appl Stat 45:516–526
- Stein ML (1999) Interpolation of spatial data: some theory for Kriging. Springer, Berlin
- Sun Y, Genton MG (2011) Functional boxplots. J Comput Grap Stat 20:313-334
- Sun Y, Genton MG (2012a) Adjusted functional boxplots for spatio-temporal data visualization and outlier detection. Environmetrics 23:54–64
- Sun Y, Genton MG (2012b) Functional median polish. J Agric Biol Environ Stat 17:354–376
- Sun Y, Genton MG, Nychka D (2012) Exact fast computation of band depth for large functional datasets: how quickly can one million curves be ranked? Stat 1:68–74