



## S-Estimators for Functional Principal Component Analysis

Graciela Boente & Matías Salibian-Barrera

**To cite this article:** Graciela Boente & Matías Salibian-Barrera (2015) S-Estimators for Functional Principal Component Analysis, Journal of the American Statistical Association, 110:511, 1100-1111, DOI: [10.1080/01621459.2014.946991](https://doi.org/10.1080/01621459.2014.946991)

**To link to this article:** <http://dx.doi.org/10.1080/01621459.2014.946991>



View supplementary material [↗](#)



Accepted author version posted online: 05 Aug 2014.  
Published online: 07 Nov 2015.



Submit your article to this journal [↗](#)



Article views: 343



View related articles [↗](#)



View Crossmark data [↗](#)

# S-Estimators for Functional Principal Component Analysis

Graciela BOENTE and Matías SALIBIAN-BARRERA

Principal component analysis is a widely used technique that provides an optimal lower-dimensional approximation to multivariate or functional datasets. These approximations can be very useful in identifying potential outliers among high-dimensional or functional observations. In this article, we propose a new class of estimators for principal components based on robust scale estimators. For a fixed dimension  $q$ , we robustly estimate the  $q$ -dimensional linear space that provides the best prediction for the data, in the sense of minimizing the sum of robust scale estimators of the coordinates of the residuals. We also study an extension to the infinite-dimensional case. Our method is consistent for elliptical random vectors, and is Fisher consistent for elliptically distributed random elements on arbitrary Hilbert spaces. Numerical experiments show that our proposal is highly competitive when compared with other methods. We illustrate our approach on a real dataset, where the robust estimator discovers atypical observations that would have been missed otherwise. Supplementary materials for this article are available online.

KEY WORDS: Functional data analysis; Robust estimation; Sparse data.

## 1. INTRODUCTION

Principal component analysis (PCA) is a widely used method to obtain a lower-dimensional approximation to multivariate data. This approximation is optimal in the sense of minimizing the mean squared loss between the original observations and the resulting approximations. Estimated principal components can be a valuable tool to explore the data visually, and are also useful to describe some characteristics of the data (e.g., directions of high variability). Thanks to the ever reducing cost of collecting data, many datasets in current applications are both large and complex, sometimes with a very high number of variables. The chance of having outliers or other type of imperfections in the data increases both with the number of observations and their dimension. Thus, detecting these outlying observations is an important step, even when robust estimates are used, either as a preprocessing step or because there is some specific interest in finding anomalous observations.

As a motivation, consider the problem of identifying days with an atypical concentration of ground level ozone ( $O_3$ ) in the air. Ground level ozone forms as a result of the reaction between sunlight, nitrogen oxide ( $NO_x$ ), and volatile organic compounds (VOC). We obtained hourly average concentration of ground level ozone at a monitoring station in Richmond, BC (a few kilometers south of the city Vancouver, BC). The data come from the Ministry of Environment of the province of British Columbia, and are available online at <http://envistaweb.env.gov.bc.ca>. We focus on the month of August for the years 2004 to 2012. Figure 1 displays the data.

Each line corresponds to the evolution of the hourly average concentration (in ppb) of ground level ozone for 1 day. A few days exceed the maximum desired level threshold of 50 ppb set by the Canadian National Ambient Air Quality Objectives, but there may also be days exhibiting an hourly pattern different from the majority of the curves.

In this article, we study robust low-dimensional approximations for high-(or infinite-) dimensional data that can be used to identify poorly fitted observations as potential outliers. The earliest and probably most immediate approach to obtain robust estimates for the principal components consists in using the eigenvalues and eigenvectors of a robust scatter estimator (Campbell 1980; Devlin, Gnanadesikan, and Kettenring 1981; Boente 1987; Naga and Antille 1990; Croux and Haesbroeck 2000). A different approach was proposed by Locantore et al. (1999) based on using the covariance matrix of the data projected onto the unit sphere. Since principal component directions are also those that provide projections with the largest variability, robust PCA estimators can alternatively be obtained as the directions that maximize a robust estimator of scale of the projected data. This approach is called “projection pursuit,” see Li and Chen (1985), Croux and Ruiz-Gazen (1996, 2005), Hubert, Rousseeuw, and Verboven (2002), and Hubert, Rousseeuw, and Vanden Branden (2005).

It is well known that, for finite-dimensional observations with finite second moments, when using mean squared errors, the best lower-dimensional approximation is given by the projections onto the linear space spanned by the eigenvectors of the covariance matrix corresponding to its largest eigenvalues. Several robust proposals exist in the literature exploiting this characterization of PCA. They amount to replacing the squared residuals with a different loss function. Liu et al. (2003) used the absolute value of the residuals, and McCoy and Tropp (2011) proposed a randomized algorithm to find an approximate solution to this  $L_1$  minimization problem. Croux et al. (2003) proposed a weighted version of this procedure that reduces the effect of high-leverage

Graciela Boente is Full Professor, Departamento de Matemáticas, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Ciudad Universitaria, Pabellón 1, Buenos Aires 1428, Argentina (E-mail: [gboente@dm.uba.ar](mailto:gboente@dm.uba.ar)). She also has a researcher position at the CONICET. Matías Salibian-Barrera is Associate Professor, Department of Statistics, University of British Columbia, 3182 Earth Sciences Building, 22007 Main Mall, Vancouver, BC, V6T 1Z4, Canada (E-mail: [matias@stat.ubc.ca](mailto:matias@stat.ubc.ca)). This research was partially supported by Grants PIP 112-201101-00339 from CONICET, PICT 0397 from ANPCYT, and w276 from the Universidad de Buenos Aires at Buenos Aires, Argentina (G. Boente) and Discovery Grant of the Natural Sciences and Engineering Research Council of Canada (M. Salibian Barrera). The authors thank the associate editor and three anonymous referees for valuable comments that led to an improved version of the original article.

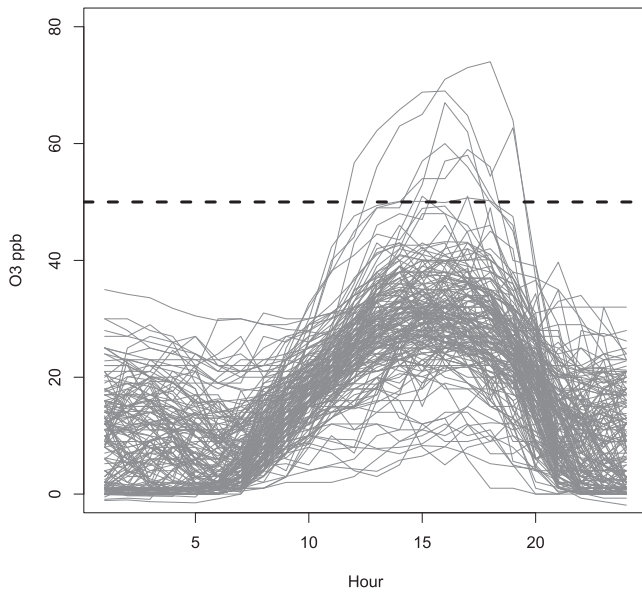


Figure 1. Hourly mean concentration (in ppb) of ground level ozone in Richmond, BC, Canada, for August of 2004 to 2012. The darker dashed horizontal line at 50 ppb is the current maximum desired level set by the Canadian National Ambient Air Quality Objectives.

points. Verboon and Heiser (1994) and De la Torre and Black (2001) used a bounded loss function applied to column-wise standardized residuals. Later, Maronna and Yohai (2008) proposed a similar loss function, but modified in such a way that the method reduces to the classical PCA when one uses a squared loss function. Maronna (2005) also considered best-estimating lower-dimensional subspaces directly, but his approach cannot be easily extended to infinite-dimensional settings because there may be infinitely many minimum eigenvalues.

There has been recent attention paid to a similar problem in the engineering and computer science literature. The main assumption in that approach is that a proportion of the observations lies on a proper lower-dimensional subspace, and that there may be a sparse amount of arbitrary additive and diffuse “noise” present. The objective is to fully recover the low-rank part of the data. Chandrasekaran et al. (2011), Candès et al. (2011), McCoy and Tropp (2011), and Xu, Caramanis, and Sanghavi (2012) studied different convex relaxations of the problem of finding an exact representation of the data matrix as the sum of a low-rank one and a sparse one. Lerman et al. (2014) and Zhang and Lerman (2014) also considered convex relaxations of this problem. The focus of these proposals is on obtaining fast algorithms, and they derive sufficient conditions to guarantee that the solution to the surrogate convex optimization problem is the lower-dimensional subspace that properly contains the “nonoutlying” points.

Our approach relies on a probabilistic model and assumes that our observations follow an elliptical distribution. We are interested in studying best lower-dimensional approximations, in the sense of minimizing the expected prediction error over the distribution of the random vector. These approximations need not fit exactly any subset of the data. Moreover, our goal is to obtain robust alternatives for estimating principal spaces in infinite-dimensional settings. We use finite (or high-)dimensional esti-

mators as a step toward achieving that purpose. Nevertheless, our proposal provides consistent estimators of the best lower-dimensional subspace when applied to multivariate data that follow an elliptical distribution, even if second moments do not exist. Furthermore, our approach is Fisher consistent for the case of infinite-dimensional observations. Few robust principal components estimates for functional data (FPCA) have been proposed in the literature. Gervini (2008) studied spherical principal components, and Hyndman and Ullah (2007) discussed a projection-pursuit approach using smoothed trajectories, but without studying their properties in detail. More recently, Sawant, Billor, and Shin (2012) adapted the BACON-PCA method to detect outliers and to provide robust estimators of the functional components, while Bali et al. (2011) proposed robust projection-pursuit FPCA estimators and showed that they are consistent to the eigenfunctions and eigenvalues of the underlying process.

The rest of the article is organized as follows. Section 2 tackles the problem of providing robust estimators for a  $q$ -dimensional approximation for Euclidean data. Section 3 discusses extending this methodology to accommodate functional data, and its use to detect outliers is described in Section 4. In Section 5 we report the results of a simulation study conducted to study the performance of the proposed procedure for functional data. The Richmond Ozone dataset is analyzed in Section 6, where the advantage of the proposed procedure to detect possible influential observations is illustrated. Finally, Section 7 provides some further discussion and recommendations. Proofs are relegated to the online supplementary materials where we also analyze the French mortality data.

## 2. S-ESTIMATORS OF THE PRINCIPAL COMPONENTS IN $\mathbb{R}^p$

Consider the problem of finding a lower-dimensional approximation to a set of observations  $\mathbf{x}_i$ ,  $1 \leq i \leq n$ , in  $\mathbb{R}^p$ . Specifically, we search for  $q < p$  vectors  $\mathbf{b}^{(l)} \in \mathbb{R}^p$ ,  $1 \leq l \leq q$ , whose spanned linear subspace provides a good approximation to the data. From now on,  $\mathbf{B} \in \mathbb{R}^{p \times q}$  stands for the matrix  $\mathbf{B} = (\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(q)})$ ,  $\mathbf{b}_j^T$  denotes the  $j$ th row of  $\mathbf{B}$ , and the subspace spanned by its columns is  $\mathcal{L}_{\mathbf{B}}$ . For a given  $\boldsymbol{\mu} \in \mathbb{R}^p$ , the corresponding “fitted values” are  $\hat{\mathbf{x}}_i = \boldsymbol{\mu} + \mathbf{B} \mathbf{a}_i$ ,  $1 \leq i \leq n$ , where  $\mathbf{a}_i \in \mathbb{R}^q$ . The principal components are defined as the minimizers, over matrices  $\mathbf{A} \in \mathbb{R}^{n \times q}$ ,  $\mathbf{B} \in \mathbb{R}^{p \times q}$ , and vectors  $\boldsymbol{\mu} \in \mathbb{R}^p$ , of

$$L_2(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}) = \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_{\mathbb{R}^p}^2 = \sum_{i=1}^n \sum_{j=1}^p r_{ij}^2, \quad (1)$$

where the  $i$ th row of the matrix  $\mathbf{A} \in \mathbb{R}^{n \times q}$  is  $\mathbf{a}_i$ ,  $r_{ij} = x_{ij} - \hat{x}_{ij}$  and  $\|\cdot\|_{\mathbb{R}^p}$  denotes the usual Euclidean norm in  $\mathbb{R}^p$ . Furthermore, this optimization problem can be solved using alternating regression iterations. Note that if we restrict  $\mathbf{B}$  to satisfy  $\mathbf{B}^T \mathbf{B} = \mathbf{I}_q$ , then the vectors  $\mathbf{a}_i$ ,  $1 \leq i \leq n$ , correspond to the scores of the sample on this basis.

Our approach is based on noting that  $L_2(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  in (1) is proportional to  $\sum_{j=1}^p s_j^2$ , where  $s_j^2$  is the sample variance of the residuals’  $j$ th coordinate:  $r_{1j}, r_{2j}, \dots, r_{nj}$ . To reduce the influence of atypical observations, we propose to use robust scale estimates instead of sample variances. Our robustly estimated

$q$ -dimensional subspace best approximating the data are defined as the linear space  $\mathcal{L}_{\mathbf{B}}$  where  $(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  minimizes

$$L_S(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}) = \sum_{j=1}^p \widehat{\sigma}_j^2, \quad (2)$$

and  $\widehat{\sigma}_j$  denotes a robust scale estimator of the residuals  $r_{ij} = x_{ij} - \widehat{x}_{ij}$ ,  $1 \leq i \leq n$ . Note that if we use the sample variance  $s_j^2$  instead of  $\widehat{\sigma}_j^2$ , then the objective function in (2) reduces to the classical one in (1).

Scale estimators measure the spread of a sample and are invariant under translations and equivariant under scale transformations (see, e.g., Maronna, Martin, and Yohai 2006). Although any robust scale estimator can be used in (2), to fix ideas we focus our presentation on  $M$ -estimators of scale (see Huber and Ronchetti 2009). As in Maronna, Martin, and Yohai (2006), let  $\rho: \mathbb{R} \rightarrow \mathbb{R}_+$  be a  $\rho$ -function, that is, an even function, nondecreasing on  $|x|$ , increasing for  $x > 0$  when  $\rho(x) < \lim_{t \rightarrow \infty} \rho(t)$  and such that  $\rho(0) = 0$ . Given residuals  $r_{ij}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}) = x_{ij} - \widehat{x}_{ij}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  with  $\widehat{x}_{ij}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}) = \mu_j + \mathbf{a}_i^T \mathbf{b}_j$ , the  $M$ -estimator of scale of the residuals  $\widehat{\sigma}_j = \widehat{\sigma}_j(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  satisfies

$$\frac{1}{n} \sum_{i=1}^n \rho_c \left( \frac{r_{ij}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})}{\widehat{\sigma}_j} \right) = b, \quad (3)$$

where  $\rho_c(u) = \rho(u/c)$ , and  $c > 0$  is a user-chosen tuning constant. When  $\rho(y) = \min(3y^2 - 3y^4 + y^6, 1)$ , (Tukey's biweight function) with  $c = 1.54764$  and  $b = 1/2$ , the estimator is Fisher consistent at the normal distribution and has breakdown point 50%. In general, if  $\|\rho\|_{\infty} = 1$ , then the breakdown point of the  $M$ -scale estimator is  $\min(b, 1 - b)$ .

We can write our estimator in a slightly more general way as follows. Let  $\pi(\mathbf{y}, \mathcal{L}_{\mathbf{B}})$  denote the orthogonal projection of  $\mathbf{y}$  onto  $\mathcal{L}_{\mathbf{B}}$ . To simplify the presentation, assume that  $\boldsymbol{\mu}$  is known. For each observation  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $1 \leq i \leq n$ , let  $\mathbf{r}_i(\mathcal{L}_{\mathbf{B}}) = \mathbf{x}_i - \boldsymbol{\mu} - \pi(\mathbf{x}_i - \boldsymbol{\mu}, \mathcal{L}_{\mathbf{B}}) = (r_{i1}(\mathcal{L}_{\mathbf{B}}), \dots, r_{ip}(\mathcal{L}_{\mathbf{B}}))^T$  denote the corresponding vector of residuals and  $\widehat{\sigma}_{j, \mathcal{L}_{\mathbf{B}}} = \widehat{\sigma}(r_{1j}(\mathcal{L}_{\mathbf{B}}), \dots, r_{nj}(\mathcal{L}_{\mathbf{B}}))$  the scale estimator of the  $j$ th coordinate of the residuals. Let  $\widehat{\Psi}_n(\mathcal{L}_{\mathbf{B}}) = \sum_{j=1}^p \widehat{\sigma}_{j, \mathcal{L}_{\mathbf{B}}}^2$ . The  $S$ -estimator of the  $q$ -dimensional principal subspace is the linear space  $\widehat{\mathcal{L}} = \widehat{\mathcal{L}}_{\mathbf{B}}$  that solves

$$\widehat{\mathcal{L}}_{\mathbf{B}} = \underset{\dim(\mathcal{L}_{\mathbf{B}})=q}{\operatorname{argmin}} \widehat{\Psi}_n(\mathcal{L}_{\mathbf{B}}). \quad (4)$$

To study the asymptotic properties of robust estimators, it is convenient to think of them as functionals of the empirical distribution of the sample (Huber and Ronchetti 2009). For example,  $M$ -scale estimators in (3) correspond to the functional  $\sigma_{\mathbf{R}}: \mathcal{D} \rightarrow \mathbb{R}_+$  defined for each distribution function  $F \in \mathcal{D}$  as the solution  $\sigma_{\mathbf{R}}(F)$  to the equation  $\int \rho_c(t/\sigma_{\mathbf{R}}(F)) dF(t) = b$ . Here,  $\mathcal{D}$  is a subset of all the univariate distributions, which contains all the empirical ones.

In what follows we will assume that  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $1 \leq i \leq n$  are independent and identically distributed random vectors with distribution  $P$ . The independence condition may be relaxed, for instance, requiring stationarity and a mixing condition or just ergodicity, since we only need the strong law of large numbers to hold to guarantee the consistency results given below. For a random vector  $\mathbf{x}$  with distribution  $P$ , denote  $F_j(\mathcal{L}_{\mathbf{B}})$  the distribu-

tion of the  $j$ th coordinate of  $\mathbf{r}(\mathcal{L}_{\mathbf{B}})$  and let  $\Psi(\mathcal{L}) = \sum_{j=1}^p \sigma_{j, \mathcal{L}}^2$ , where  $\sigma_{j, \mathcal{L}} = \sigma_{\mathbf{R}}(F_j(\mathcal{L}_{\mathbf{B}}))$ . The functional  $\mathcal{L}(P)$  corresponding to the  $S$ -estimators defined in (4) is the linear space of dimension  $q$  that satisfies

$$\mathcal{L}(P) = \underset{\dim(\mathcal{L})=q}{\operatorname{argmin}} \Psi(\mathcal{L}). \quad (5)$$

Recall that a random vector is said to have a spherical distribution if its distribution is invariant under orthogonal transformations. In particular, the characteristic function of a spherically distributed  $\mathbf{x} \in \mathbb{R}^p$  is of the form  $\varphi_{\mathbf{x}}(\mathbf{t}) = \phi(\mathbf{t}^T \mathbf{t})$  for  $\mathbf{t} \in \mathbb{R}^p$ , where  $\phi: \mathbb{R} \rightarrow \mathbb{R}$  is the generator of the characteristic function. We write  $\mathbf{x} \sim \mathcal{S}_p(\phi)$ . For a  $p \times p$  matrix  $\mathbf{B}$  and a vector  $\boldsymbol{\mu} \in \mathbb{R}^p$ , the distribution of  $\mathbf{x} = \mathbf{B}\mathbf{z} + \boldsymbol{\mu}$  when  $\mathbf{z} \sim \mathcal{S}_p(\phi)$  is called elliptical,  $\mathcal{E}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \phi)$ , where  $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}^T$ . The following proposition, whose proof is relegated to the online supplementary materials, shows that the solution to (5) is the desired linear space.

**Proposition 2.1.** Let  $\mathbf{x} \sim \mathcal{E}_p(\mathbf{0}, \boldsymbol{\Sigma}, \phi)$  be an elliptically distributed random vector with  $\boldsymbol{\Sigma} = \boldsymbol{\beta}\boldsymbol{\Lambda}\boldsymbol{\beta}^T$ ,  $\boldsymbol{\Lambda} = \operatorname{diag}(\lambda_1, \dots, \lambda_p)$ ,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ , where  $\boldsymbol{\beta}$  is an orthonormal matrix with columns  $\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(p)}$ . Assume that  $\lambda_q > \lambda_{q+1}$ . Then, the linear space  $\mathcal{L}_q$  spanned by  $\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(q)}$  is the unique solution of (5).

As mentioned before, this approach can also be used with any robust scale estimator. For example, we can define  $\tau$ -estimators by considering the  $\tau$ -best lower-dimensional approximations, given by the minimizers of  $L_{\tau}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}) = \sum_{j=1}^p \widehat{\sigma}_j^2 \sum_{i=1}^n \rho_1(r_{ij}(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})/\widehat{\sigma}_j)$ , where  $\widehat{\sigma}_j = \widehat{\sigma}_j(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  is an  $M$ -scale estimator computed as in (3) with a  $\rho$ -function  $\rho$  such that  $\rho \leq \rho_1$ . Note that if an iterative procedure is used to solve (4), the scale estimators  $\widehat{\sigma}_j$  need to be updated at each step of the algorithm.

Consistency of projection-pursuit principal component estimators for random vectors was derived in Cui, He, and Ng (2003) requiring uniform convergence over the unit ball of the projected data scale estimators to the scale functional. This condition was generalized in Bali et al. (2011) to the functional case. A natural extension for  $q > 1$  is

$$\sup_{\dim(\mathcal{L})=q} |\widehat{\Psi}_n(\mathcal{L}) - \Psi(\mathcal{L})| \xrightarrow{\text{a.s.}} 0. \quad (6)$$

Note that this condition is easily verified when using a robust scale functional with finite-dimensional random vectors since the Stiefel manifold  $\mathcal{V}_{p \times q} = \{\mathbf{B} \in \mathbb{R}^{p \times q} : \mathbf{B}^T \mathbf{B} = \mathbf{I}_q\}$  is a compact set. Furthermore, the following proposition shows that this condition is sufficient to obtain consistency of the  $S$ -estimators in (4).

**Proposition 2.2.** Assume that  $\mathcal{L}(P)$  is unique and that (6) holds. Then, the estimators  $\widehat{\mathcal{L}} = \widehat{\mathcal{L}}_{\mathbf{B}}$  obtained minimizing  $\widehat{\Psi}_n(\mathcal{L})$  in (4) over linear spaces  $\mathcal{L}$  of dimension  $q$ , are consistent to the linear space  $\mathcal{L}(P)$  defined in (5). In other words, with probability one,  $\pi(\mathbf{x}, \widehat{\mathcal{L}})$  converges to  $\pi(\mathbf{x}, \mathcal{L}(P))$ , for almost all  $\mathbf{x}$ .

## 2.1 Algorithm for $S$ -Estimators

The optimization problem defining our estimator is generally nonconvex, and typically difficult to solve. In this section, we show that first-order conditions for a critical point of the objective function in (4) naturally suggest an iterative reweighted



least-square algorithm. Once such iterations are available, a standard strategy used in the statistical literature to compute this type of estimators (e.g., Rousseeuw and van Driessen 1999; Maronna 2005; Salibian-Barrera and Yohai 2006) is to use a large number of random initial points, and select the best visited local minimum as the estimator.

Note that although  $S$ -scale estimators are only defined implicitly, explicit first-order conditions can be obtained differentiating both sides of (3). More specifically, let  $\hat{\sigma}_j$ ,  $j = 1, \dots, p$  be an  $M$ -estimator of scale of the residuals  $x_{ij} - \hat{x}_{ij}$ ,  $i = 1, \dots, n$ . In other words,  $\hat{\sigma}_j$  satisfies  $(1/n) \sum_{i=1}^n \rho((x_{ij} - \mu_j - \mathbf{a}_i^T \mathbf{b}_j) / \hat{\sigma}_j) = b$ , where we have absorbed the constant  $c$  into the loss function  $\rho$ . The derivatives with respect to  $\mathbf{a}_i$ ,  $i = 1, \dots, n$  are given by

$$\begin{aligned} \frac{\partial}{\partial \mathbf{a}_i} \left( \sum_{j=1}^p \hat{\sigma}_j^2 \right) &= \sum_{j=1}^p 2 \hat{\sigma}_j \frac{\partial \hat{\sigma}_j}{\partial \mathbf{a}_i} \\ &= -2 \sum_{j=1}^p \hat{\sigma}_j h_j^{-1} \rho' \left( \frac{r_{ij}}{\hat{\sigma}_j} \right) \mathbf{b}_j, \\ i &= 1, \dots, n, \end{aligned}$$

where  $h_j = \sum_{i=1}^n \rho'(r_{ij}/\hat{\sigma}_j) r_{ij}/\hat{\sigma}_j$ . Similarly, the other first-order conditions are

$$\begin{aligned} \frac{\partial}{\partial \mathbf{b}_s} \left( \sum_{j=1}^p \hat{\sigma}_j^2 \right) &= \sum_{j=1}^p 2 \hat{\sigma}_j \frac{\partial \hat{\sigma}_j}{\partial \mathbf{b}_s} \\ &= -2 \hat{\sigma}_s h_s^{-1} \sum_{i=1}^n \rho' \left( \frac{r_{is}}{\hat{\sigma}_s} \right) \mathbf{a}_i, \\ s &= 1, \dots, p \\ \frac{\partial}{\partial \mu_\ell} \left( \sum_{j=1}^p \hat{\sigma}_j^2 \right) &= \sum_{j=1}^p 2 \hat{\sigma}_j \frac{\partial \hat{\sigma}_j}{\partial \mu_\ell} \\ &= -2 \hat{\sigma}_\ell h_\ell^{-1} \sum_{i=1}^n \rho' \left( \frac{r_{i\ell}}{\hat{\sigma}_\ell} \right), \\ \ell &= 1, \dots, p. \end{aligned}$$

Setting these to zero, we obtain a system of equations that can be reexpressed as reweighted least-square problems as follows: let  $w_{ij} = \hat{\sigma}_j h_j^{-1} \rho'(r_{ij}/\hat{\sigma}_j)$ , then we need to solve

$$\begin{aligned} \sum_{j=1}^p w_{ij} (x_{ij} - \mu_j) \mathbf{b}_j &= \left( \sum_{j=1}^p w_{ij} \mathbf{b}_j \mathbf{b}_j^T \right) \mathbf{a}_i, \quad 1 \leq i \leq n, \\ \sum_{i=1}^n w_{ij} (x_{ij} - \mu_j) \mathbf{a}_i &= \left( \sum_{i=1}^n w_{ij} \mathbf{a}_i \mathbf{a}_i^T \right) \mathbf{b}_j, \quad 1 \leq j \leq p, \\ \sum_{i=1}^n w_{ij} (x_{ij} - \mathbf{a}_i^T \mathbf{b}_j) &= \sum_{i=1}^n w_{ij} \mu_j, \quad 1 \leq j \leq p. \end{aligned}$$

This formulation suggests the usual iterative reweighted least-square (IRWLS) algorithm. Given initial estimates  $\mathbf{b}_j^{(0)}$ ,  $1 \leq j \leq p$ , and  $\boldsymbol{\mu}^{(0)}$  compute the scores  $\mathbf{a}_i^{(0)}$ ,  $i = 1, \dots, n$ , the weights  $w_{ij}^{(0)}$  and obtain updated values for  $\mathbf{a}_i^{(1)}$ ,  $\mathbf{b}_j^{(1)}$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq p$ , and  $\boldsymbol{\mu}^{(1)}$ . We repeat these steps until the objective function changes less than a chosen tolerance value. The best  $q$ -dimensional linear space approximation is spanned by

$\{\hat{\mathbf{b}}^{(1)}, \dots, \hat{\mathbf{b}}^{(q)}\}$ , the final values obtained above. For interpretation purposes, we orthogonalize the set  $\{\hat{\mathbf{b}}^{(1)}, \dots, \hat{\mathbf{b}}^{(q)}\}$  and compute the scores  $\hat{\mathbf{a}}_i$  as the corresponding orthogonal projections.

For the initial location vector  $\boldsymbol{\mu}^{(0)}$ , we use the  $L^1$ -median, and adapt the strategy of Rousseeuw and van Driessen (1999) to select initial values for  $\mathbf{B}$  and  $\mathbf{A}$ . More specifically, we generate  $N_1$  random starts for the matrix  $\mathbf{B}$ , which are orthogonalized, each of them leading to an initial matrix  $\mathbf{B}^{(0)}$ . The columns of the matrix  $\mathbf{A}$  are the scores of each observation on the basis given by the  $q$  columns of  $\mathbf{B}^{(0)}$ . For each of these initial values, we run  $N_2$  IRWLS iterations, or until a tolerance level is achieved. The initial values giving the best objective function after  $N_2$  iterations are then iterated until convergence. This algorithm depends on the number of random starts  $N_1$ , the desired tolerance for sequential change in the objective function, and the number of iterations  $N_2$  that is applied to each random candidate. In our experiments, we used a tolerance of  $10^{-6}$  and found that using  $N_1 = 50$  random starts and  $N_2 = 50$  partial IRWLS iterations for each of them was typically sufficient to find a good solution to (4), which is in line with the results of Maronna (2005).

An implementation of this algorithm in R is publicly available online from <http://www.stat.ubc.ca/~matias/soft.html>. Although a formal computational complexity analysis of this algorithm is beyond the scope of this article, our numerical experiments reported in Section 5 show that the algorithm works very well. We tested the speed of our R code using these settings on an Intel i7 CPU (3.5GHz) machine running Windows 7. In Table 1, we report the average time in CPU minutes over 10 random samples for different combinations of the sample size ( $n$ ), number of variables ( $p$ ), and dimension of the subspace ( $q$ ). Note that these times could be improved notably if the algorithm was implemented in C or a language with faster linear algebra operations.

## 2.2 Choosing the Dimension of the Approximating Subspace

In some cases, the desired dimension of the linear subspace providing an approximation to the data is either known or chosen in advance (e.g., for visualization purposes). In many applications, however, this dimension is selected based on the resulting “proportion of unexplained variability.”

Proposition 2.1 shows that for  $\mathbf{x} \sim \mathcal{E}_p(\mathbf{0}, \boldsymbol{\Sigma}, \mathbb{E})$ , the functional  $\Psi(\mathcal{L})$  is minimized, when  $\mathcal{L} = \mathcal{L}_q$  the subspace spanned by the first  $q$  eigenvectors of the scatter matrix and  $\Psi(\mathcal{L}_q) = \sum_{j=q+1}^p \lambda_j$ . Note that for  $q = 0$ , we have  $\Psi(\mathcal{L}_0) = \sum_{j=1}^p \lambda_j = \text{tr}(\boldsymbol{\Sigma}) = \sum_{j=1}^p \sigma_{j,0}^2$ , where  $\sigma_{j,0} = \sigma_{\mathbf{R}}(F_{j,0})$  with  $F_{j,0}$  the distribution of  $r_j(\boldsymbol{\mu}) = x_j - \mu_j$ . Thus, the proportion of unexplained variability can be defined as  $u_q = \Psi(\mathcal{L}_q)/\Psi(\mathcal{L}_0)$  and an estimator of  $u_q$  is given by  $\hat{u}_q = \hat{\Psi}_n(\hat{\mathcal{L}}_q)/\hat{\Psi}_n(\hat{\mathcal{L}}_0)$ , where  $\hat{\mathcal{L}}_q$  is defined in (4) and  $\hat{\mathcal{L}}_0$  corresponds to minimizing  $\hat{\Psi}_n(\mathcal{L}_0) = \sum_{j=1}^p \hat{\sigma}_{j,\mathcal{L}_0}^2$  with  $\hat{\sigma}_{j,\mathcal{L}_0} = \hat{\sigma}(r_{1j}(\boldsymbol{\mu}), \dots, r_{nj}(\boldsymbol{\mu}))$  the scale estimator of the  $j$ th coordinate of the residuals  $\mathbf{r}_i(\boldsymbol{\mu}) = \mathbf{x}_i - \boldsymbol{\mu}$ . Proposition 2.2 can be used to show the consistency of  $\hat{u}_q$  to  $u_q$ .

To avoid the high computational cost of solving (4) for different values of  $q$ , we adapt the strategy of Maronna (2005). Let  $u_{\max}$  be the maximum allowed proportion of unexplained variability, and a maximum dimension  $q_{\max}$  of the approximating subspace. We look for the smallest  $q_0$  such that  $q_0 \leq q_{\max}$

Table 1. Average timing of the IRWLS algorithm (in CPU minutes)

<i>n</i>	<i>p</i> = 50			<i>p</i> = 100			<i>p</i> = 200			<i>p</i> = 500		
	<i>q</i> = 1	<i>q</i> = 2	<i>q</i> = 5	<i>q</i> = 1	<i>q</i> = 2	<i>q</i> = 5	<i>q</i> = 1	<i>q</i> = 2	<i>q</i> = 5	<i>q</i> = 1	<i>q</i> = 2	<i>q</i> = 5
50	4.2	4.3	3.9	8.7	8.9	8.1	17.8	17.9	16.8	53.5	53.8	53.3
100	5.0	4.9	4.8	9.9	10.0	8.9	20.5	22.6	23.6	69.0	67.2	70.8
200	5.6	6.0	5.8	11.5	12.2	10.8	25.8	28.1	25.0	97.6	108.5	116.3

and  $\widehat{u}_{q_0} \leq u_{\max}$ . We first verify that  $\widehat{u}_{q_{\max}} \leq u_{\max}$  otherwise the problem cannot be solved and we need to modify our goals. The procedure starts with  $q_1 = 1$ . If  $\widehat{u}_1 \leq u_{\max}$  we are done. Otherwise, assume that after  $j$  steps, we have  $\widehat{u}_{q_j} \geq u_{\max}$ , where  $q_j = j$  dimension used in step  $j$ . Let  $\widehat{\mu}^{(q_j)}$  be the estimated center and  $\widehat{\mathbf{B}}_{q_j} \in \mathbb{R}^{p \times q_j}$  the orthonormal basis of the best  $q_j$ -dimensional subspace, with columns  $\widehat{b}_{q_j}^{(1)}, \dots, \widehat{b}_{q_j}^{(q_j)}$ . As before, let  $\widehat{\mathbf{A}}_{q_j} \in \mathbb{R}^{n \times q_j}$  be the matrix of scores. Let  $q_{j+1} = q_j + 1$  and define the matrices  $\mathbf{B} = (\widehat{\mathbf{B}}_{q_j}, \boldsymbol{\beta}) \in \mathbb{R}^{p \times q_{j+1}}$ , with  $\boldsymbol{\beta} \in \mathbb{R}^p$ , and  $\mathbf{A} = (\widehat{\mathbf{A}}_{q_j}, \boldsymbol{\alpha}) \in \mathbb{R}^{n \times q_{j+1}}$  with  $\boldsymbol{\alpha} \in \mathbb{R}^n$ . Let  $\mathbf{b}_1, \dots, \mathbf{b}_{q_j+1}$  and  $\mathbf{a}_1, \dots, \mathbf{a}_n$  denote the columns of  $\mathbf{B}$  and the rows of  $\mathbf{A}$ , respectively. We construct our predictions as  $\widehat{x}_{i\ell}^{(q_{j+1})} = \widehat{\mu}_{\ell}^{(q_j)} + \mathbf{a}_i^T \mathbf{b}_{\ell}$ , and note that the residuals satisfy  $r_{i\ell}^{(q_{j+1})} = r_{i\ell}^{(q_j)} - \alpha_i \beta_{\ell}$ . Our problem is now to minimize  $L_S(\mathbf{A}, \mathbf{B}, \widehat{\mu}^{(q_j)})$  over  $\boldsymbol{\beta}, \boldsymbol{\alpha}$  such that  $\widehat{\mathbf{B}}_{q_j}^T \boldsymbol{\beta} = \mathbf{0}$ , with  $L_S(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  given in (2). A system of equations analogous to that described in Section 2.1 can be derived to formulate an iterative reweighted least-square algorithm. Once the optimal  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$  are found, we optimize  $L_S(\mathbf{A}, \mathbf{B}, \boldsymbol{\mu})$  over  $\boldsymbol{\mu}$  to obtain  $\widehat{\mu}^{(q_{j+1})}$ . This approach is much faster than solving (4) for  $q = q_{j+1}$ . Note that  $\widetilde{u}_{q_{j+1}} = \widehat{\Psi}_n(\widehat{\mathcal{L}}_{q_{j+1}}) / \widehat{\Psi}_n(\widehat{\mathcal{L}}_0)$  is typically larger than  $\widehat{u}_{q_{j+1}}$ , so that if  $\widetilde{u}_{q_{j+1}} \leq u_{\max}$ , we select  $q = q_{j+1}$ , and otherwise increase  $j$  and continue.

3. S-ESTIMATORS IN THE FUNCTIONAL SETTING

In this section, we discuss extensions of the estimators defined in Section 2 to accommodate functional data. The most common situation corresponds to the case when the observations correspond to realizations of a stochastic process  $X \in L^2(\mathcal{I})$  with  $\mathcal{I}$  an interval of the real line, which can be assumed to be  $\mathcal{I} = [0, 1]$ . A more general setup that can accommodate applications where observations are images, for example, is to consider realizations of a random element on a separable Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  and norm  $\| \cdot \|_{\mathcal{H}}$ . Note that principal components for functional data (defined via the Karhunen–Loève decomposition of the covariance function of the process  $X$ ) also have the property of providing best lower-dimensional approximations,

in the  $L^2$  sense. Recently, a stochastic best lower-dimensional approximation for elliptically distributed random elements on separable Hilbert spaces, such as those considered when dealing with multivariate data, was obtained by Boente, Salibián-Barrera, and Tyler (2014). This optimality property does not require second moment conditions.

However, even in the simplest situation when  $X \in L^2([0, 1])$ , one rarely observes entire curves. The functional datum for replication  $i$  usually corresponds to a finite set of discrete values  $x_{i1}, \dots, x_{im_i}$  with  $x_{ij} = X_i(t_{ij})$ ,  $1 \leq j \leq m_i$ . Depending on the characteristics of the grid of points  $t_{ij}$  where observations were obtained, one can employ different strategies to analyze these data.

The easiest situation is when observations were made at common design points. In this case, we have  $p = m_1 = m_i$  and  $t_{ij} = \tau_j$ , for all  $1 \leq i \leq n$  and  $1 \leq j \leq p$ . Defining  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$ , a purely multivariate approach can be used as in Section 2 to obtain a  $q$ -dimensional linear space  $\widehat{\mathcal{L}}$  spanned by orthonormal vectors  $\widehat{\mathbf{b}}^{(1)}, \dots, \widehat{\mathbf{b}}^{(q)}$ . An associated basis in  $L^2([0, 1])$  can be defined as  $\widehat{\phi}_{\ell}(\tau_j) = a_{\ell} \widehat{b}_{\ell j}$ , for  $1 \leq \ell \leq q$ ,  $1 \leq j \leq p$ , where  $a_{\ell}$  is a constant to ensure that  $\| \widehat{\phi}_{\ell} \|_{L^2} = 1$  and  $\widehat{\mathbf{b}}^{(\ell)} = (b_{\ell 1}, \dots, b_{\ell p})^T$ . Smoothing over the observed data points, one can recover the complete trajectory. This approach provides a consistent estimator for the best approximating linear space and the corresponding “fitted trajectories”  $\pi(X_i, \widehat{\mathcal{L}})$ ,  $1 \leq i \leq n$ .

In many cases, however, trajectories are observed at different design points  $t_{ij}$ ,  $1 \leq j \leq m_i$ ,  $1 \leq i \leq n$ . In what follows, we will assume that as the sample size  $n$  increases, so does the number of points where each trajectory is observed and that, in the limit, these points cover the interval  $[0, 1]$ . Our approach consists of using a sequence of finite-dimensional functional spaces, which increases with the sample size. The basic idea is to identify each observed point in  $\mathcal{H}$  with the vector formed by its coordinates on a finite-dimensional basis that increases with the sample size. The procedure of Section 2 can be applied to these vectors to obtain a  $q$ -dimensional approximating subspace, which can then be mapped back onto  $\mathcal{H}$ .

Table 2. Mean prediction errors over 500 replications for Model 1

Method	$\epsilon_1 = \epsilon_2 = 0.00$		$\epsilon_1 = 0.10$			$\epsilon_1 = 0.20$			
	Clean	Out	Clean	Out	Clean	Out	Clean	Out	Clean
True	1.266	26.930	1.138	269.316	1.264	53.780	1.013	269.685	1.265
LS	1.246	18.961	5.065	193.372	5.679	37.429	5.682	187.461	7.104
S (3)	1.253	26.922	1.126	269.245	1.252	53.425	1.081	268.453	1.361
S (1.5)	1.308	26.872	1.270	268.937	1.417	53.241	1.464	267.400	1.850
PP	1.335	26.536	1.335	265.791	1.486	51.845	1.559	260.972	1.972

Table 3. Mean prediction errors over 500 replications for Model 2

Method	$\epsilon_1 = \epsilon_2 = 0.00$		$\epsilon_1 = 0.10$			$\epsilon_1 = 0.20$			
	Clean	Out	Clean	Out	Clean	Out	Clean	Out	Clean
True	1.359	10.063	1.222	100.589	1.358	20.054	1.087	100.598	1.358
LS	1.339	1.597	4.032	19.528	4.512	1.840	4.482	9.505	5.610
S (3)	1.346	9.839	1.380	99.230	1.541	12.427	2.357	69.919	3.035
S (1.5)	1.401	9.638	2.047	97.207	2.296	17.916	2.891	90.648	3.645
PP	1.428	8.922	1.427	90.696	1.589	14.865	1.618	76.535	2.039

More specifically, let  $\{\delta_i\}_{i \geq 1}$  be an orthonormal basis of  $\mathcal{H}$  and, for each  $n \geq 1$ , let  $\mathcal{H}_{p_n}$  be the linear space spanned by  $\delta_1, \dots, \delta_{p_n}$ . To simplify the notation, we write  $p = p_n$ . Let  $x_{ij} = \langle X_i, \delta_j \rangle_{\mathcal{H}}$  be the coefficient of the  $i$ th trajectory on the  $j$ th element of the basis,  $1 \leq j \leq p$ , and form the  $p$ -dimensional vector  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$ . When,  $\mathcal{H} = L^2([0, 1])$ , the inner products  $\langle X_i, \delta_j \rangle_{\mathcal{H}}$  can be numerically computed using a Riemann sum over the design points for the  $i$ th trajectory  $\{t_{ij}\}_{1 \leq j \leq p}$ . We apply the procedure described in Section 2 to the multivariate observations  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  to obtain a  $q$ -dimensional linear space  $\hat{\mathcal{L}}$  spanned by orthonormal vectors  $\hat{\mathbf{b}}^{(1)}, \dots, \hat{\mathbf{b}}^{(q)}$  and the corresponding “predicted values”  $\hat{\mathbf{x}}_i = \hat{\boldsymbol{\mu}} + \sum_{\ell=1}^q \hat{a}_{i\ell} \hat{\mathbf{b}}^{(\ell)}$ , with  $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \dots, \hat{\mu}_p)^T$ . It is now easy to find the corresponding approximation in the original space  $\mathcal{H}$ . The location parameter is  $\hat{\mu}_{\mathcal{H}} = \sum_{j=1}^p \hat{\mu}_j \delta_j$ , and the associated  $q$ -dimensional basis in  $\mathcal{H}$  is  $\hat{\phi}_{\ell} = \sum_{j=1}^p \hat{b}_{\ell j} \delta_j / \|\sum_{j=1}^p \hat{b}_{\ell j} \delta_j\|_{\mathcal{H}}$ , for  $1 \leq \ell \leq q$ . Furthermore, the “fitted values” in  $\mathcal{H}$  are  $\hat{X}_i = \hat{\mu}_{\mathcal{H}} + \sum_{\ell=1}^q \hat{a}_{i\ell} \hat{\phi}_{\ell}$ . Moreover, since  $\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_{\mathbb{R}^p} \simeq \|X_i - \hat{X}_i\|_{\mathcal{H}}$ , we can also use squared residual norms to detect atypical observations.

For a proof of the Fisher consistency of this approach, we refer the reader to Section S.1 of the online supplementary materials.

### 3.1 Algorithm for Functional Data

In this section, we give details on how to compute our S-estimators for functional principal components. The basic idea consists of applying the algorithm of Section 2.1 to the coordinates of the observed data on a sufficiently rich orthonormal basis of the Hilbert space, and then transforming back the result to the original variables.

To fix ideas, consider the case where the data consist of functions  $X_i$ ,  $1 \leq i \leq n$ , observed at points  $t_1, \dots, t_m$ . We approximate the  $L^2$  inner product with a Riemann sum over the grid of points:  $\langle \alpha, \beta \rangle_{\mathcal{H}} = \int \alpha(t)\beta(t) dt \approx \sum_{\ell=2}^m \alpha(t_{\ell})\beta(t_{\ell})(t_{\ell} - t_{\ell-1})$ . Let  $v_1, \dots, v_p$  be a  $B$ -spline basis. We orthonormalize  $v_1, \dots, v_p$  using the approximated inner product to obtain or-

thonormal elements  $\delta_1, \dots, \delta_p$ . Let  $\mathbf{\Delta} \in \mathbb{R}^{m \times p}$  be the matrix of the functions  $\delta_j$  evaluated at the points  $t_i$ :  $\mathbf{\Delta} = (\delta_1, \delta_2, \dots, \delta_p)$ , where  $\delta_j = (\delta_j(t_1), \delta_j(t_2), \dots, \delta_j(t_m))^T$ . Then, if  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is the matrix of observed trajectories (one in each row), the coordinates of each  $X_i$  on each element  $\delta_j$  of the spline basis is denoted as  $\tilde{\mathbf{x}}_{i,j} = \sum_{\ell=2}^m X_i(t_{\ell})\delta_j(t_{\ell})(t_{\ell} - t_{\ell-1}) \approx \langle X_i, \delta_j \rangle_{\mathcal{H}}$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq p$ . We now apply the algorithm given in Section 2.1 to the “data” matrix  $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times p}$  of the coordinates of our observations on the  $B$ -spline basis. We obtain the center vector  $\tilde{\boldsymbol{\mu}} \in \mathbb{R}^p$ , an orthonormal basis  $\tilde{\mathbf{B}} \in \mathbb{R}^{p \times q}$  of the best  $q$ -dimensional subspace, and the matrix of scores  $\tilde{\mathbf{A}} \in \mathbb{R}^{n \times q}$ . The matrix  $\tilde{\mathbf{X}} = \mathbb{I}_n \tilde{\boldsymbol{\mu}}^T + \tilde{\mathbf{A}} \tilde{\mathbf{B}}^T$  provides the  $q$ -dimensional approximation to our functional data written in the  $B$ -splines basis. Finally, we express our solution in the original variables  $\hat{\mathbf{X}} = \tilde{\mathbf{X}} \mathbf{\Delta}^T$ . Note that  $\hat{\mathbf{X}} = \mathbb{I}_n (\mathbf{\Delta} \tilde{\boldsymbol{\mu}})^T + \tilde{\mathbf{A}} (\mathbf{\Delta} \tilde{\mathbf{B}})^T$ . In other words,  $\mathbf{\Delta} \tilde{\boldsymbol{\mu}} \in \mathbb{R}^m$  is the vector of the center function  $\hat{\mu}_{\mathcal{H}}$  evaluated at the points  $t_1, \dots, t_m$ , and  $\mathbf{\Delta} \tilde{\mathbf{B}} \in \mathbb{R}^{m \times q}$  is the matrix of  $q$  orthonormal functions  $\hat{\phi}_{\ell}$  spanning the best lower approximation space in  $\mathcal{H}$ , evaluated on the same points.

## 4. OUTLIER DETECTION

An important use of robust estimators for multivariate data is the detection of potential outliers; see, for example, Rousseeuw and Van Zomeren (1990), Becker and Gather (2001), Pison and van Aelst (2004), and Hardin and Rocke (2005). Unfortunately, these approaches do not extend naturally to the functional case.

Alternatively, one can consider the PCA residuals as indicators of outlyingness. Given a sample  $\mathbf{x}_1, \dots, \mathbf{x}_n$  in  $\mathbb{R}^p$  and the estimated subspace  $\hat{\mathcal{L}} = \mathcal{L}_{\hat{\mathbf{B}}}$  in (4), one can construct the corresponding “best  $q$ -dimensional” approximations  $\hat{\mathbf{x}}_i = \hat{\boldsymbol{\mu}} + \pi(\mathbf{x}_i - \hat{\boldsymbol{\mu}}, \mathcal{L}_{\hat{\mathbf{B}}}) = \hat{\boldsymbol{\mu}} + \hat{\mathbf{B}} \hat{\mathbf{B}}^T (\mathbf{x}_i - \hat{\boldsymbol{\mu}})$ ,  $1 \leq i \leq n$ . We expect outlying or otherwise atypical observations to be poorly fitted and thus to have a relatively large residual  $R_i = \|\mathbf{r}_i(\mathcal{L}_{\hat{\mathbf{B}}})\|_{\mathbb{R}^p} = \|(\mathbf{I} - \hat{\mathbf{B}} \hat{\mathbf{B}}^T)(\mathbf{x}_i - \hat{\boldsymbol{\mu}})\|_{\mathbb{R}^p}$ ,  $1 \leq i \leq n$ . Exploring the norm of these residuals sometimes provides sufficient information to

Table 4. Mean prediction errors over 500 replications for Model 3

Method	$\epsilon = 0.00$		$\epsilon = 0.10$			$\epsilon = 0.20$			
	Clean	Out	Clean	Out	Clean	Out	Clean	Out	Clean
True	0.304	4.411	0.274	44.163	0.304	8.842	0.243	44.088	0.304
LS	0.285	2.074	0.660	18.457	0.736	5.599	0.711	27.363	0.893
S (3)	0.301	4.412	0.269	44.148	0.299	8.846	0.237	44.113	0.297
S (1.5)	0.354	4.465	0.318	44.674	0.354	8.931	0.284	44.535	0.355
PP	0.385	4.439	0.355	44.397	0.394	8.913	0.321	44.430	0.402

Table 5. Average sensitivity and specificity over 500 random samples following Model 1

$\epsilon_1$	LS	PP	S (3)	S (1.5)	HDR	BAG	LRT	DTR	DWE	HU
Sensitivity										
0.10	0.914	1.000	1.000	0.998	0.155	0.597	0.305	1.000	1.000	1.000
0.20	0.295	0.835	0.856	0.833	0.074	0.224	0.018	1.000	1.000	1.000
Specificity										
0.00	0.982	0.982	0.981	0.982	0.986	0.983	0.978	0.802	0.802	0.782
0.10	0.999	0.997	0.996	0.997	0.999	0.982	1.000	0.839	0.839	0.792
0.20	1.000	1.000	1.000	1.000	0.999	0.989	1.000	0.897	0.897	0.808

detect abnormal points in the data. It is worth noticing that the distribution of the residuals squared norm  $R_i^2$  is unknown, but typically skewed to the right because they are bounded by 0 from below. Following the approach of Hubert and Vandervieren (2008), we propose to flag an observation as atypical if its squared residual norm exceeds the upper whisker of a skewed-adjusted boxplot.

Another way to use principal components to look for potential outliers considers the scores of each point on the estimated principal eigenvectors. The solution to (4) provides an estimated basis  $\hat{\mathbf{b}}^{(j)}$ ,  $1 \leq j \leq q$  (the columns of  $\hat{\mathbf{B}}$ ) for the optimal  $q$ -dimensional linear space spanned by the first  $q$  eigenvectors, but the  $\hat{\mathbf{b}}^{(j)}$ 's themselves need not be estimates of the principal directions. However, we can use an approach similar to "projection pursuit" to sequentially search for vectors in  $\hat{\mathcal{L}}_{\hat{\mathbf{B}}}$  that maximize a robust scale estimate of the corresponding projections of the data. Specifically, for each  $\boldsymbol{\gamma} \in \hat{\mathcal{L}}_{\hat{\mathbf{B}}}$ , let  $F_n[\boldsymbol{\gamma}]$  be the empirical distribution of the projected observations  $\boldsymbol{\gamma}^T \mathbf{x}_1, \dots, \boldsymbol{\gamma}^T \mathbf{x}_n$ , and  $\sigma_R(F_n[\boldsymbol{\gamma}])$  be the corresponding scale estimator. The estimated first principal direction is obtained maximizing  $\sigma_R(F_n[\boldsymbol{\gamma}])$  over unitary vectors in  $\hat{\mathcal{L}}_{\hat{\mathbf{B}}}$ . Subsequent principal directions are similarly computed with the additional condition of being orthogonal to the previous ones. The scores of each observation on the estimated principal directions can be used to screen for atypical data points.

Both of these last two approaches have natural counterparts for functional data and can be used with the estimators defined in Section 3. Hyndman and Shang (2010) defined two detection rules based on the scores of a robust two-dimensional fit and compared them with a residuals-based PCA procedure introduced by Hyndman and Ullah (2007). Our simulation study in Section 5 includes these methods as well those based on functional depth proposed by Febrero, Galeano, and Gonzalez-Manteiga (2007, 2008).

As in the finite-dimensional case, to find potential outliers one may consider looking for curves  $X_i$  that are poorly predicted by the  $S$ -estimator using the squared prediction errors  $R_{i,\mathcal{H}}^2 = \|X_i - \hat{X}_i\|_{\mathcal{H}}^2$ ,  $i = 1, \dots, n$ . As in the finite-dimensional case, the distribution of these prediction residuals is unknown and difficult to estimate. Hyndman and Ullah (2007) proposed to use a normal approximation to the residual squared norm, which they called the integrated squared error, to define a threshold. Our approach is more data analytic and does not depend on the underlying distribution of the process even if we always have in mind that the uncontaminated process has an elliptical distribution. For that reason, we mimic the proposal given in the finite-dimensional case and to decide whether an observation may be flagged as a potential outlier, we used the adjusted boxplot of Hubert and Vandervieren (2008) on the residuals  $R_{i,\mathcal{H}}^2$ , identifying as an atypical observation a value exceeding the upper whisker of the adjusted boxplot. We use this approach in the example and in our simulation study discussed below.

## 5. SIMULATION

In this section, we present the results of a simulation study performed to investigate the finite-sample properties of our robust sieve proposal. In all cases, we generated 500 samples of size  $n = 70$  where each trajectory was observed at  $m = 100$  equidistant points in the interval  $[0, 1]$ . We used a cubic  $B$ -spline basis of dimension  $p = 50$ , which is sufficiently rich to represent the data well. This choice represents a realistic situation where the sample size is similar to the dimension of the problem. Other reasonable choices for the dimension of the spline basis (even with  $n < p$ ) yielded very similar results and lead to the same conclusions in our numerical experiments.

Table 6. Average sensitivity and specificity over 500 random samples following Model 2

$\epsilon_1$	LS	PP	S (3)	S (1.5)	HDR	BAG	LRT	DTR	DWE	HU
Sensitivity										
0.10	0.178	0.996	0.979	0.915	0.135	0.774	0.059	0.350	0.353	1.000
0.20	0.020	0.708	0.637	0.474	0.053	0.079	0.005	0.239	0.239	1.000
Specificity										
0.00	0.980	0.980	0.980	0.980	0.986	0.982	0.978	0.803	0.803	0.782
0.10	0.996	0.997	0.997	0.997	0.997	0.958	1.000	0.817	0.817	0.774
0.20	0.994	1.000	0.997	1.000	0.994	0.988	0.999	0.815	0.815	0.770



Table 7. Average sensitivity and specificity over 500 random samples following Model 3

$\epsilon_1$	LS	PP	S (3)	S (1.5)	HDR	BAG	LRT	DTR	DWE	HU
Sensitivity										
0.10	0.936	1.000	1.000	1.000	0.148	0.489	0.124	0.982	0.988	1.000
0.20	0.603	0.848	0.850	0.848	0.071	0.418	0.063	0.922	0.977	1.000
Specificity										
0.00	0.987	0.986	0.987	0.987	0.986	0.983	0.990	0.804	0.804	0.849
0.10	0.998	0.997	0.998	0.998	0.999	0.988	1.000	0.838	0.837	0.869
0.20	1.000	1.000	1.000	1.000	0.999	0.991	1.000	0.863	0.886	0.896

## 5.1 Simulation Settings

The following three different models constructed from finite- and infinite-range processes were used to generate the data. In two of them we included a relatively small proportion of measurement errors, as is usual in many applications.

*Model 1.* This model corresponds to the case where most of the curves follow a smooth trajectory, but some of them may display sudden vertical jumps at a few time points. In this setup, the noncontaminated observations  $X_i \sim X$ ,  $1 \leq i \leq n$ , with  $X(t_s) \sim 10 + \mu(t_s) + \xi_1 \phi_1(t_s) + \xi_2 \phi_2(t_s) + z_s$ ,  $s = 1, \dots, 100$ , where the additive errors  $z_s$  are iid  $N(0, 1)$ , the scores  $\xi_1 \sim N(0, 25/4)$ ,  $\xi_2 \sim N(0, 1/4)$ ,  $\xi_1$  and  $\xi_2$  are independent and independent of  $z_s$ . The mean function is  $\mu(t) = 5 + 10 \sin(4\pi t) \exp(-2t) + 5 \sin(\pi t/3) + 2 \cos(\pi t/2)$  and  $\phi_1(t) = \sqrt{2} \cos(2\pi t)$  and  $\phi_2(t) = \sqrt{2} \sin(2\pi t)$  correspond to the Fourier basis.

We also generated contaminated trajectories  $X_i^{(c)}$  as realizations of the process  $X^{(c)}$  defined by  $X^{(c)}(t_s) = X(t_s) + V Y(t_s)$ ,  $s = 1, \dots, 100$ , where  $V \sim \text{Bi}(1, \epsilon_1)$  is independent of  $X$  and  $Y$ ,  $Y(t_s) = W_s \tilde{z}_s$  with  $W_s \sim \text{Bi}(1, \epsilon_2)$ ,  $\tilde{z}_s \sim N(\mu^{(c)}, 0.01)$ ,  $W_s$  and  $\tilde{z}_s$  are all independent. In other words, a trajectory is contaminated with probability  $\epsilon_1$ , and at any point  $t_s$  the contaminated function is shifted with probability  $\epsilon_2$ . The shift is random but tightly distributed around the constant  $\mu^{(c)} = 30$ . Samples without outliers correspond to  $\epsilon_1 = 0$ . To investigate the influence of different outlier configurations of our estimator, we considered the settings:  $\epsilon_1 = 0.10$  and  $\epsilon_1 = 0.20$ , with  $\epsilon_2 = 0.30$  in both cases.

*Model 2* This situation corresponds to a similar case as in Model 1, but with some curves starting on a different trajectory that joins smoothly with the one that most curves follow. In this case, noncontaminated observations  $X_i \sim X$  were generated as  $X(t_s) \sim 150 - 2\mu(t_s) + \xi_1 \phi_1(t_s) + \xi_2 \phi_2(t_s) + z_s$ ,  $s = 1, \dots, 100$ , where  $z_s$ ,  $\xi_1$ ,  $\xi_2$ ,  $\mu$ ,  $\phi_1$ , and  $\phi_2$  are as in the previous model. However, contaminated trajectories are only perturbed in a specific part of their range. The atypical observations satisfy  $X_i^{(c)} \sim X^{(c)}$ , where  $X^{(c)}(t_s) = X(t_s) + V Y(t_s)$  for  $t_s < 0.4$  and  $X^{(c)}(t_s) = X(t_s)$  for  $t_s \geq 0.4$ , where  $V \sim \text{Bi}(1, \epsilon_1)$  is independent of  $X$  and  $Y$ ,  $Y(t_s) = W_s \tilde{z}_s$  with  $W_s \sim \text{Bi}(1, \epsilon_2)$ ,  $\tilde{z}_s \sim N(\mu^{(c)}(t_s), 0.01)$ , with  $\mu^{(c)}(t_s) = -5 - 2\mu(t_s)$ , and  $W_s$  and  $\tilde{z}_s$  are all independent. In this model, we used  $\epsilon_1 = 0.10$  and  $\epsilon_1 = 0.20$ , and in both cases we set  $\epsilon_2 = 0.90$ .

*Model 3.* This setting corresponds to functions that follow an infinite-rank stochastic process. Contamination is present in terms of short, sudden vertical shifts. Curves were generated from a Gaussian process with covariance kernel  $\gamma_X(s, t) = 10 \min(s, t)$ . The eigenfunctions of the covariance operator equal  $\phi_j(t) = \sqrt{2} \sin((2j-1)(\pi/2)t)$ ,  $j \geq 1$ , with associated eigenvalues  $\lambda_j = 10(2/[d(2j-1)\pi])^2$ . As in Sawant, Billor, and Shin (2012), the contaminated observations  $X_i^{(c)}$  are defined as  $X_i^{(c)}(s) = X_i(s) + V_i D_i M \mathbb{I}_{[T_i < s < T_i + \ell]}$ , where  $V_i \sim \text{Bi}(1, \epsilon)$ ,  $\mathbb{P}(D_i = 1) = \mathbb{P}(D_i = -1) = 1/2$ ,  $T_i \sim \mathcal{U}(0, 1 - \ell)$ ,  $\ell < 1/2$ , and  $V_i$ ,  $X_i$ ,  $D_i$ , and  $T_i$  are independent. We choose  $\ell = 1/15$ ,  $M = 30$ , and  $\epsilon = 0.1$  and  $0.2$ .

## 5.2 The Estimators

We computed the classical principal components estimator (LS) as well as the robust one defined in (2), using an  $M$ -scale estimator, with function  $\rho_c$  in Tukey's bisquare family with tuning constants  $c = 1.54764$  and  $b = 0.50$ . We also considered the choice  $c = 3.0$  and  $b = 0.2426$ , which we expect to yield more efficiency. The robust estimators are labeled as S (1.5) and S (3) in the tables. As mentioned in Section 2.1, after obtaining the robust  $q$ -dimensional linear space, we orthonormalize its basis and compute the scores  $\hat{\mathbf{a}}_i$  as the corresponding orthogonal projections. We also computed the sieve projection-pursuit approach proposed in Bali et al. (2011), which is called "PP" in our tables. For comparison purposes, we have also calculated the mean squared prediction errors obtained with the true best  $q$ -dimensional linear space for uncontaminated data. This benchmark is indicated as "True" in all tables. Since trajectories following Models 1 and 2 were generated using a two-dimensional scatter operator (i.e., the underlying process had only two nonzero eigenvalues) plus measurement errors, we used  $q = 1$  with our estimator. For Model 3, we used  $q = 4$ , which results in 95% of explained variance.

## 5.3 Simulation Results

To summarize the results of our simulation study, for each replication we consider mean squared prediction errors in the original space, that is, based on  $\|X_i - \hat{X}_i\|_{\mathcal{H}}^2$ . The conclusions that can be reached using the finite-dimensional residuals squared prediction error  $\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_{\mathbb{R}^p}^2$  are the same as those discussed below, and hence are not reported here. Tables 2 to 4 report the average mean squared error for outlying and nonoutlying trajectories separately, as a way to quantify how the

procedures fit the bulk of the data. More specifically, let  $\gamma_i = 1$  when  $X_i$  is an outlier and  $\gamma_i = 0$  otherwise, then

$$\text{PE}_{\mathcal{H},\text{OUT}} = \frac{1}{n} \sum_{i=1}^n \gamma_i \|X_i - \hat{X}_i\|_{\mathcal{H}}^2$$

and

$$\text{PE}_{\mathcal{H},\text{CLEAN}} = \frac{1}{n} \sum_{i=1}^n (1 - \gamma_i) \|X_i - \hat{X}_i\|_{\mathcal{H}}^2. \quad (7)$$

Note that the total prediction error equals  $\text{PE}_{\mathcal{H}} = (1/n) \sum_{i=1}^n \|X_i - \hat{X}_i\|_{\mathcal{H}}^2 = \text{PE}_{\mathcal{H},\text{OUT}} + \text{PE}_{\mathcal{H},\text{CLEAN}}$ . We also report the mean PE over contaminated and clean trajectories separately:

$$\overline{\text{PE}}_{\mathcal{H},\text{OUT}} = \frac{\sum_{i=1}^n \gamma_i \|X_i - \hat{X}_i\|_{\mathcal{H}}^2}{\sum_{i=1}^n \gamma_i}$$

and

$$\overline{\text{PE}}_{\mathcal{H},\text{CLEAN}} = \frac{\sum_{i=1}^n (1 - \gamma_i) \|X_i - \hat{X}_i\|_{\mathcal{H}}^2}{\sum_{i=1}^n (1 - \gamma_i)}.$$

We also compute the prediction squared errors of the actual best lower-dimensional predictions  $\hat{X}_i^0$ , obtained with the first  $q$  true eigenfunctions (recall that we used  $q = 1$  in Models 1 and 2, and  $q = 4$  in Model 3). The results for this “estimator” are tabulated in the row labeled “True.” The averages over the 500 replications of  $\text{PE}_{\mathcal{H},\text{OUT}}$ ,  $\text{PE}_{\mathcal{H},\text{CLEAN}}$ ,  $\overline{\text{PE}}_{\mathcal{H},\text{OUT}}$ , and  $\overline{\text{PE}}_{\mathcal{H},\text{CLEAN}}$  are labeled “Out,” “Clean,” “Out,” and “Clean,” respectively.

As expected, when no outliers are present all procedures are comparable, with a small loss for the robust procedures. The  $S$ -estimator with  $c = 3$  had the second smallest mean squared prediction error, after the LS. When samples were contaminated, the classical procedure based on least squares tries to compromise between outlying and nonoutlying trajectories and this is reflected on the values of  $\text{PE}_{\mathcal{H},\text{OUT}}$  and  $\text{PE}_{\mathcal{H},\text{CLEAN}}$  in (7) and also on the average prediction error of the contaminated and noncontaminated trajectories  $\overline{\text{PE}}_{\mathcal{H},\text{OUT}}$  and  $\overline{\text{PE}}_{\mathcal{H},\text{CLEAN}}$ . With contaminated samples, the  $S$ -estimator had the best performance overall. Its mean squared prediction was closest to the “True” one, and it also provided better fits to the noncontaminated samples (and worse predictions for the contaminated trajectories). This last observation can be seen comparing the columns labeled “Out” and “Clean.” The only case when the sieves projection-pursuit estimator performed slightly better than the  $S$ -estimator is for Model 2 with  $\epsilon_1 = 0.20$  and  $\epsilon_2 = 0.90$  (see Table 3). The

advantage of the  $S$ -estimator was more notable in all the other cases of Model 1, Model 2, and Model 3.

We also compared the performance of different outlier detection methods for functional data. As described in Section 4, we used the squared prediction errors  $R_{i,\mathcal{H}}^2 = \|X_i - \hat{X}_i\|_{\mathcal{H}}^2$ ,  $i = 1, \dots, n$ , to find curves  $X_i$  that are poorly predicted. Those with squared prediction errors exceeding the upper whisker of the adjusted boxplot will be flagged as outliers. We used the same approach with predictors  $\hat{X}_i$  obtained using the other estimators mentioned before.

In addition, we included other outlier-detection methods for functional data that appeared in the literature. We considered the functional high-density region and the functional bagplots of Hyndman and Shang (2010) with a 99% coverage, denoted as HDR and BAG, respectively, as well as the integrated squared error method defined in Hyndman and Ullah (2007), denoted as HU. The first two methods are based on the scores of a two-dimensional robust projection-pursuit fit. To keep the comparison fair, for HU we chose a  $q$ -dimensional robust fit with  $q = 1$  under Models 1 and 2 and  $q = 4$  under Model 3. Furthermore, we also compared our detection rule with the proposals based on a likelihood-ratio-type statistic given in Febrero, Galeano, and Gonzalez-Manteiga (2007) and on the modal depth, using both trimmed and weighted bootstrap estimates for the threshold as proposed in Febrero, Galeano, and Gonzalez-Manteiga (2008). These methods are denoted as LRT, DTR, and DWE, respectively. These detection rules are implemented in the R package *rainbow*.

For each model and each outlier detection method, in Tables 5 to 7 we report the average sensitivity and specificity over the 500 samples. Sensitivity is the proportion of actual outliers that are correctly flagged as such, while specificity is the proportion of nonoutlying curves correctly identified as not atypical. An ideal method will simultaneously maintain high sensitivity and specificity.

For Model 1, we note that DRT, DWE, and HU identify too many curves as outliers (resulting in a high sensitivity but low specificity). On the other hand, LRT, HDR, and BAG consistently miss most of the outliers (low sensitivity), as does LS when the proportion of outliers is 20%. Using prediction residuals based on  $S$ - and the projection-pursuit estimators offers the

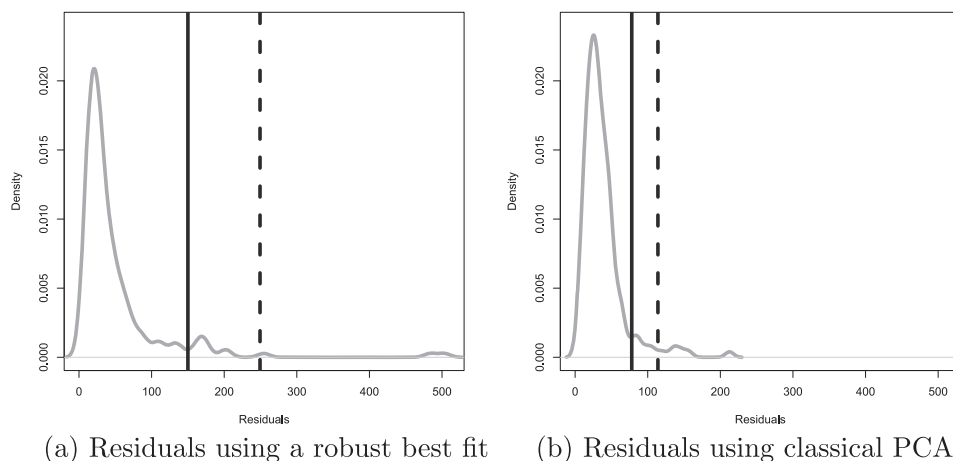


Figure 2. Estimated density of the squared prediction errors with (a) the  $S$ -estimator and (b) the classical one. The dashed line corresponds to the threshold suggested by `adjbox()` while the solid one indicates the beginning of a relatively heavy tail.

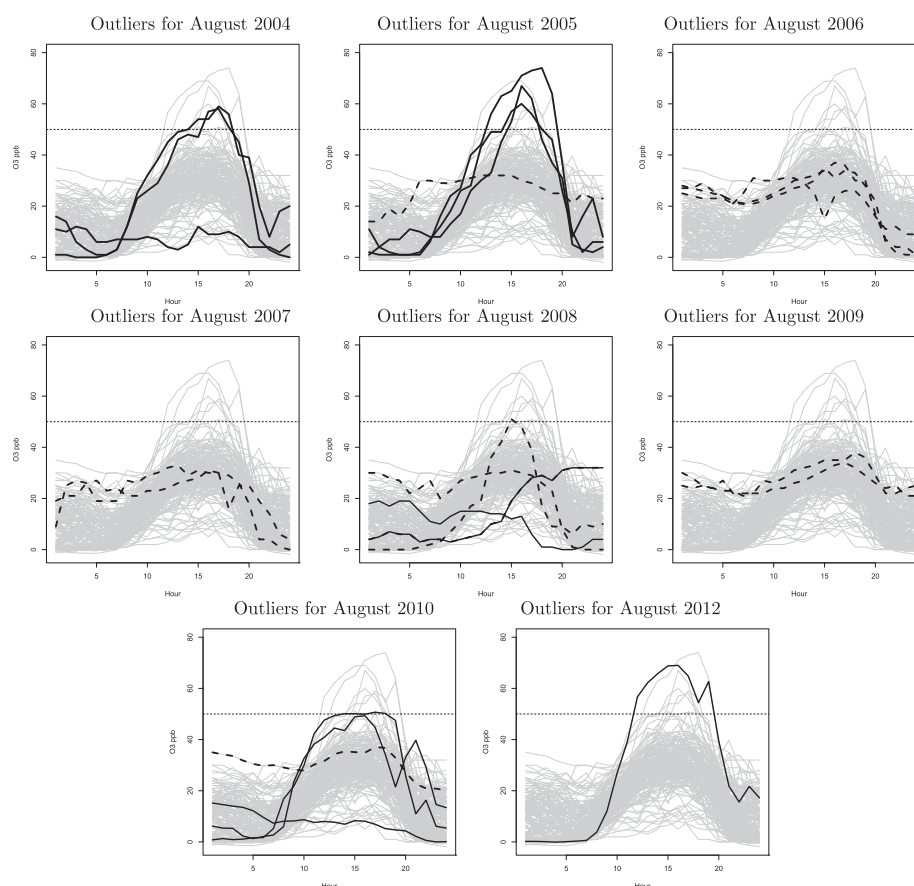


Figure 3. Hourly mean concentration (in ppb) of ground level ozone in Richmond, BC, Canada. Thin gray lines show all the available data. Solid lines correspond to potential outliers identified by the robust estimator, dashed lines to those found by a classical analysis.

best overall performance. When the data follow Model 2, LS, HDR, LRT, DTR, and DWE fail to detect most of the outliers, as does BAG for  $\epsilon = 0.20$ . Again, HU flags too many curves as outlying. The relatively low specificity of DTR and DWE (and to some extent BAG) seems to indicate that the few observations flagged as outliers are not the truly atypical ones. Once again the approach based on  $S$ - and projection-pursuit estimators works best. Note that although the  $S(1.5)$  appears to miss around half of the outliers for  $\epsilon_1 = 0.20$ , those flagged as atypical are correctly identified. The results for Model 3 are very similar to those for Model 1. Overall, for the three scenarios considered here, the clear best method to detect functional outliers is to use the squared prediction residuals based on a robust principal components estimator.

## 6. EXAMPLE: GROUND LEVEL OZONE CONCENTRATIONS

These data contain hourly average measurements of ground level ozone ( $O_3$ ) concentration from a monitoring station in Richmond, BC, Canada. Ozone at ground level is a serious air pollutant and its presence typically peaks in summer months. We focus on the month of August, and obtained data for the years 2004 to 2012. We have 176 days with hourly average  $O_3$  measurements. Our purpose is to identify days in which the temporal pattern of  $O_3$  concentration appears different from the others. Based on the strong pattern observed in the data, we consider one-dimensional approximations. We use an  $S$ -estimator

with tuning constant  $c = 3$  applying the approach described in Section 3 with a cubic  $B$ -spline basis of dimension  $p = 10$ . To find potentially outlying curves, we use as threshold the upper whisker of the adjusted boxplot of Hubert and Vandervieren (2008) applied to the squared prediction errors using the LS and  $S$ -estimators. Figure 2 contains the estimated density of the  $L^2$  norm of the residuals for each of the 176 curves when we compute predictions using our  $S$ -estimators (panel (a)) and the classical LS ones (panel (b)). The dashed line in Figure 2 corresponds to the threshold suggested by the adjusted boxplot. While there are a few extreme outliers at the right tail of each plot, both plots also show a relatively heavy tail that suggests the presence of moderate outliers. The solid line indicates approximately the beginning of this heavy tail, and is the cut-off used in our analysis.

To make the visualization of the results easier, each panel in Figure 3 shows the observations detected as outliers in 1 year, both by the robust estimator (solid lines) and the classical approach (dashed lines). The thin gray lines in the background show all the available observations, and are included as a visual reference, while the light dashed horizontal line at 50 ppb is the current maximum recommended level. We see that the robust fit identifies as outliers all of the days with relatively high peaks of  $O_3$  concentration, but also some days that exhibit a “flat” profile.

Since ground level ozone is produced by the reaction between sunlight and other compounds in the air, we use temperature data to verify whether the potential outliers identified above

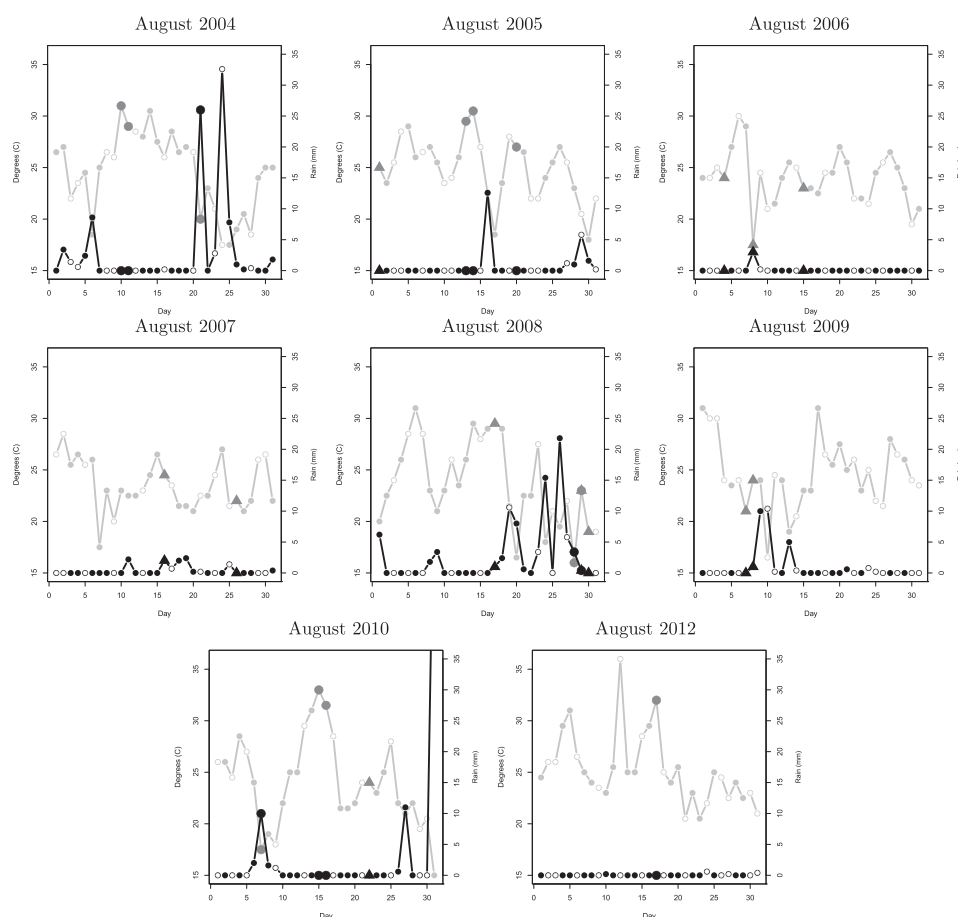


Figure 4. Maximum daily temperature profile (in black) and rain levels (in gray) for the month of August. Solid circles indicate atypical days found by the robust approach, triangles correspond to the classical method.

correspond to atypical days. Figure 4 shows maximum daily temperature for the months of August between 2004 and 2012 together with the daily amount of rain. Days for which  $O_3$  data are not available are indicated with white circles. A day identified as having an atypical  $O_3$  profile by the robust fit is marked with a large solid circle. Potential outliers identified by the classical approach are indicated with a solid triangle. We see that the outliers identified by the robust fit correspond to days with either a very high or low temperature. Furthermore, outlying days with a “flat”  $O_3$  profile are those with a low maximum temperature, while days with a sharp  $O_3$  peak correspond to particularly hot days. On the other hand, days flagged as possible outliers by LS generally do not show any pattern with respect to temperature. This analysis shows that the robust method is able to identify potential outliers that correspond to extreme values of an unobserved but closely associated meteorological variable (temperature). In other words, the robust method is able to uncover outliers that correspond to actual atypical days.

## 7. CONCLUDING REMARKS

In this article, we propose a robust estimator for the subspace spanned by the first  $q$  principal components. We show that our method is consistent and can be used in general settings, including functional data applications. In this case, our method works well when the observations can be well represented in

a sufficiently rich but arbitrary basis. Moreover, the resulting robust predictions can be used to detect atypical observations in the data. This is confirmed in our simulation study, where this outlier detection method compares very favorably to other proposals in the literature. Our estimators are defined via a non-convex optimization problem, which is difficult to solve. As it is done for similar problems arising in other contexts (robust linear regression and multivariate location and scatter estimators, e.g.), we use first-order conditions to derive an iterative reweighted least-square-type algorithm. Extensive numerical experiments show that this algorithm provides estimators with good statistical properties. It would be interesting, but beyond the scope of this work, to study whether a convex relaxation of the optimization problem (2) can provide a more scalable algorithm with comparable robustness and statistical properties.

## SUPPLEMENTARY MATERIALS

The supplementary material has three sections. Section 1 contains a discussion on the Fisher-consistency of the Sieves-approach for  $S$ -estimators for functional principal components. Section 2 includes the analysis of the French mortality dataset. Finally, in Section 3 the proofs of Propositions 2.1 and 2.2 and of the Fisher-consistency is given.

[Received June 2013. Revised July 2014.]



## REFERENCES

- Bali, L., Boente, G., Tyler, D., and Wang, J. L. (2011), "Robust Functional Principal Components: A Projection-Pursuit Approach," *The Annals of Statistics*, 39, 2852–2882. [1101,1102,1107]
- Becker, C., and Gather, U. (2001), "The Largest Nonidentifiable Outliers: A Comparison of Multivariate Simultaneous Outliers Identification Rules," *Computational Statistics and Data Analysis*, 36, 119–127. [1105]
- Boente, G. (1987), "Asymptotic Theory for Robust Principal Components," *Journal of Multivariate Analysis*, 21, 67–78. [1100]
- Boente, G., Salibian-Barrera, M., and Tyler, D. (2014), "A Characterization of Elliptical Distributions and Some Optimality Properties of Principal Components for Functional Data," *Journal of Multivariate Analysis*, 131, 254–264. [1104]
- Campbell, N. A. (1980), "Robust Procedures in Multivariate Analysis I: Robust Covariance Estimation," *Applied Statistics*, 29, 231–237. [1100]
- Candès, E. J., Li, X., Ma, Y., and Wright, J. (2011), "Robust Principal Component Analysis?" *Journal of the ACM*, 58, 1–37. [1101]
- Chandrasekaran, V., Sanghavi, S., Parrilo, P., and Willsky, A. S. (2011), "Rank-Sparsity Incoherence for Matrix Decomposition," *SIAM Journal of Optimization*, 21, 572–596. [1101]
- Croux, C., Filzmoser, P., Pison, G., and Rousseeuw, P. J. (2003), "Fitting Multiplicative Models by Robust Alternating Regressions," *Statistics and Computing*, 13, 23–36. [1100]
- Croux, C., and Haesbroeck, G. (2000), "Principal Component Analysis Based on Robust Estimators of the Covariance or Correlation Matrix: Influence Functions and Efficiencies," *Biometrika*, 87, 603–618. [1100]
- Croux, C., and Ruiz-Gazen, A. (1996), "A Fast Algorithm for Robust Principal Components Based on Projection Pursuit," in *Compstat: Proceedings in Computational Statistics*, ed. A. Prat, Heidelberg: Physica-Verlag, pp. 211–216. [1100]
- (2005), "High-Breakdown Estimators for Principal Components: The Projection-Pursuit Approach Revisited," *Journal of Multivariate Analysis*, 95, 206–226. [1100]
- Cui, H., He, X., and Ng, K. W. (2003), "Asymptotic Distribution of Principal Components Based on Robust Dispersions," *Biometrika*, 90, 953–966. [1102]
- De la Torre, F., and Black, M. J. (2001), "Robust Principal Components Analysis for Computer Vision," in *Proceedings of the 8th International Conference on Computer Vision*, 1, pp. 362–369. [1101]
- Devlin, S. J., Gnanadesikan, R., and Kettenring, J. R. (1981), "Robust Estimation of Dispersion Matrices and Principal Components," *Journal of the American Statistical Association*, 76, 354–362. [1100]
- Febrero, M., Galeano, P., and Gonzalez-Manteiga, W. (2007), "A Functional Analysis of NO<sub>x</sub> Levels: Location and Scale Estimation and Outlier Detection," *Computational Statistics*, 22, 411–427. [1106,1108]
- (2008), "Outlier Detection in Functional Data by Depth Measures, With Application to Identify Abnormal No<sub>x</sub> Levels," *Environmetrics*, 19, 331–345. [1106,1108]
- Gervini, D. (2008), "Robust Functional Estimation Using the Spatial Median and Spherical Principal Components," *Biometrika*, 95, 587–600. [1101]
- Hardin, J., and Rocke, D. (2005), "The Distribution of Robust Distances," *Journal of Computational and Graphical Statistics*, 14, 1–19. [1105]
- Huber, P. J., and Ronchetti, E. M. (2009), *Robust Statistics* (2nd ed.), New York: Wiley. [1102]
- Hubert, M., Rousseeuw, P. J., and Vanden Branden, K. (2005), "ROBPCA: A New Approach to Robust Principal Component Analysis," *Technometrics*, 47, 64–79. [1100]
- Hubert, M., Rousseeuw, P. J., and Verboven, S. (2002), "A Fast Method for Robust Principal Components With Applications to Chemometrics," *Chemometrics and Intelligent Laboratory Systems*, 60, 101–111. [1100]
- Hubert, M., and Vandervieren, E. (2008), "An Adjusted Boxplot for Skewed Distributions," *Computational Statistics and Data Analysis*, 52, 5186–5201. [1106,1109]
- Hyndman, R. J., and Ullah, S. (2007), "Robust Forecasting of mortality and Fertility Rates: A Functional Data Approach," *Computational Statistics and Data Analysis*, 51, 4942–4956. [1101,1106,1108]
- Hyndman, R. J., and Shang, H. L. (2010), "Rainbow Plots, Bagplots, and Boxplots for Functional Data," *Journal of Computational and Graphical Statistics*, 19, 29–45. [1106,1108]
- Lerman, G., McCoy, M., Tropp, J. A., and Zhang, T. (2014), "Robust Computation of Linear Models, or How to Find a Needle in a Haystack," *Foundations of Computational Mathematics*, 15, 363–410. [1101]
- Li, G., and Chen, Z. (1985), "Projection Pursuit Approach to Robust Dispersion Matrices and Principal Components: Primary Theory and Monte Carlo," *Journal of the American Statistical Association*, 80, 759–766. [1100]
- Liu, L., Hawkins, D., Ghosh, S., and Young, S. (2003), "Robust Singular Value Decomposition Analysis of Microarray Data," in *Proceedings of the National Academy of Sciences*, 100, pp. 13167–13172. [1100]
- Locantore, N., Marron, J. S., Simpson, D. G., Tripoli, N., Zhang, J. T., and Cohen, K. L. (1999), "Robust Principal Components for Functional Data," *Test*, 8, 1–28. [1100]
- Maronna, R. (2005), "Principal Components and Orthogonal Regression Based on Robust Scales," *Technometrics*, 47, 264–273. [1101,1103]
- Maronna, R., Martin, R. D., and Yohai, V. (2006), *Robust Statistics: Theory and Methods*, Chichester, UK: Wiley. [1102]
- Maronna, R., and Yohai, V. (2008), "Robust Lower-Rank Approximation of Data Matrices With Element-Wise Contamination," *Technometrics*, 50, 295–304. [1101]
- McCoy, M., and Tropp, J. A. (2011), "Two Proposals for Robust PCA Using Semidefinite Programming," *Electronic Journal of Statistics*, 5, 1123–1160. [1100,1101]
- Naga, R., and Antille, G. (1990), "Stability of Robust and Non-Robust Principal Component Analysis," *Computational Statistics and Data Analysis*, 10, 169–174. [1100]
- Pison, G., and van Aelst, S. (2004), "Diagnostic Plots for Robust Multivariate Methods," *Journal of Computational and Graphical Statistics*, 13, 1–20. [1105]
- Rousseeuw, P. J., and van Driessen, K. (1999), "A Fast Algorithm for the Minimum Covariance Determinant Estimator," *Technometrics*, 41, 212–223. [1103]
- Rousseeuw, P. J., and Van Zomeren, B. C. (1990), "Unmasking Multivariate Outliers and Leverage Points," *Journal of the American Statistical Association*, 85, 633–651. [1105]
- Salibian-Barrera, M., and Yohai, V. J. (2006), "A Fast Algorithm for S-Regression Estimates," *Journal of Computational and Graphical Statistics*, 15, 414–427. [1103]
- Sawant, P., Billor, N., and Shin, H. (2012), "Functional Outlier Detection With Robust Functional Principal Component Analysis," *Computational Statistics*, 27, 83–102. [1101,1107]
- Verboon, P., and Heiser, W. J. (1994), "Resistant Lower-Rank Approximation of Matrices by Iterative Majorization," *Computational Statistics and Data Analysis*, 18, 457–467. [1101]
- Xu, H., Caramanis, C., and Sanghavi, S. (2012), "Robust PCA via Outlier Pursuit," *IEEE Transactions on Information Theory*, 58, 3047–3064. [1101]
- Zhang, T., and Lerman, G. (2014), "A Novel M-Estimator for Robust PCA," *Journal of Machine Learning Research*, 15, 749–808. [1101]