

Exploratory Functional Data Analysis

Zhuo Qu¹, Wenlin Dai², Carolina Euan¹, Ying Sun¹, and Marc G. Genton¹

September 3, 2023

Abstract

With the advance of technology, functional data are being recorded more frequently, whether over time or in different locations. Traditionally, functional data were assumed to be defined on common and finite coordinate grids. However, real-world functional data often exhibit irregular coordinate grids or multiple components. To adapt to the demands of practical applications, researchers have developed visualization tools, outlier detection techniques, and clustering/classification methods that can handle more general types of functional data. This paper offers a comprehensive overview of recent exploratory functional data analysis (EFDA). It begins by introducing fundamental statistical concepts, such as mean and covariance functions, as well as robust statistics, such as the median and quantiles in multivariate functional data. Then, the paper delves into the evolution of visualization methods, such as the rainbow plot, and various adaptations of the functional boxplot. These modified versions of the functional boxplot are designed to accommodate the complexities of general functional data. In addition to visualization tools, the paper also reviews outlier detection technologies, which are commonly integrated with visualization methods to identify anomalous patterns within data. Moreover, the application of clustering and classification techniques tailored for functional data is explored. In closing, future directions for EFDA are briefly discussed.

Keywords: Classification, Clustering, Data visualization, Exploratory data analysis, Functional data, Multivariate functional data, Outlier detection

¹Statistics Program, King Abdullah University of Science and Technology, Thuwal 23955-6900, Saudi Arabia. E-mail: zhuo.qu@kaust.edu.sa, carolina.euancampos@kaust.edu.sa, ying.sun@kaust.edu.sa, marc.genton@kaust.edu.sa

This research was supported by the King Abdullah University of Science and Technology (KAUST).

²Center for Applied Statistics, Institute of Statistics and Big Data, Renmin University of China, Beijing 100872, China. E-mail: wenlin.dai@ruc.edu.cn

1 Introduction

Exploratory data analysis (EDA) ([Tukey 1977](#), [Martinez et al. 2017](#)) serves as the primary step in data analysis because it explores intuitively the basic properties of the underlying dataset and provides diagnostics for statistical modeling. [Tukey \(1977\)](#) contrasts EDA with confirmatory data analysis (CDA) ([Tukey 1980](#)), an area of data analysis that is mostly concerned with the techniques of statistical hypothesis testing, confidence intervals, estimation, etc. Overall, EDA can be categorized into data visualization and data mining. The data visualization tools include but are not limited to the scatter plot, the histogram, the boxplot ([Tukey 1977](#)), and the quantile-quantile plot, whereas the data mining techniques include, without limitation to the dimensionality reduction, data clustering/classification and smoothing.

When the data object changes from univariate (multivariate) data to a real (multivariate) function of an index, such as the time, wavelength, or location index, we call the new data object univariate (multivariate) functional data ([Ramsay & Dalzell 1991](#), [Ramsay & Silverman 2005](#)). Common real-life examples of univariate functional data include raw cell-cycle gene expression curves ([Zhao et al. 2004](#)), data from a longitudinal study of the relative diameter and relative height of trees ([López-Pintado & Romo 2009](#)), longitudinal height data for teenagers ([López-Pintado & Romo 2009](#)), petroleum level curves in an oil refinery ([Ramsay et al. 2009](#)), and daily temperature curves ([Sun & Genton 2011](#)), whereas real-life examples of multivariate functional data include longitudinal hip and knee angle curves for children ([Ramsay & Silverman 2005](#)), daily temperature curves derived from data collected by sensors located at different altitudes ([Berrendero et al. 2011](#)), coordinates of handwriting data ([López-Pintado et al. 2014](#)), hurricane trajectories ([Yao et al. 2005](#), [Harris](#)

et al. 2021), individual growth velocity curves for different body parts (Carroll et al. 2021), and the joint curves of stunted growth and prevalence of low-birth weight in 77 countries (Qu & Genton 2022).

Mathematically, functional data are considered as a realization of a stochastic process taking values in a Hilbert space. Each subject from one realization of the above process is assumed to be independent and to have a continuous sample path. Practically, we could never observe a function entirely over the whole domain; instead, we take records at certain fixed or random discrete points, which are either the same or differ between subjects. According to the sampling scheme, functional data can be classified as follows: 1) fully observed functions without noise at an arbitrarily dense grid (see the smoothed daily Canadian temperature curve); 2) densely sampled functions with noisy measurements (dense design, see the original daily joint Canadian temperature and precipitation curves); and 3) sparsely sampled functions with noisy measurements (sparse design in Qu & Genton 2022, see the univariate CD4 data and bivariate hurricane data). To demonstrate the application of exploratory analysis methods in diverse sampling schemes, we will use the following three dataset representatives: 1) univariate sparse CD4 cell count data from the R package *refund* (see Figure 1 (a)), 2) bivariate sparse hurricane trajectory data (downloaded [online](#)) (see Figure 1 (b)), and 3) bivariate dense Canadian daily temperature and precipitation curves from the R package *fda* (see Figure 1 (c)).

Functional data can be regarded as a natural extension of a vector from the finite dimension to the infinite dimension. However, with the continuing development of data collection techniques, functional observations present themselves more frequently. Hence, functional data analysis (FDA) (Ramsay et al. 2009) includes both an intrinsic and an applied interest. The intrinsically infinite dimension of functional data poses challenges for the existing visu-

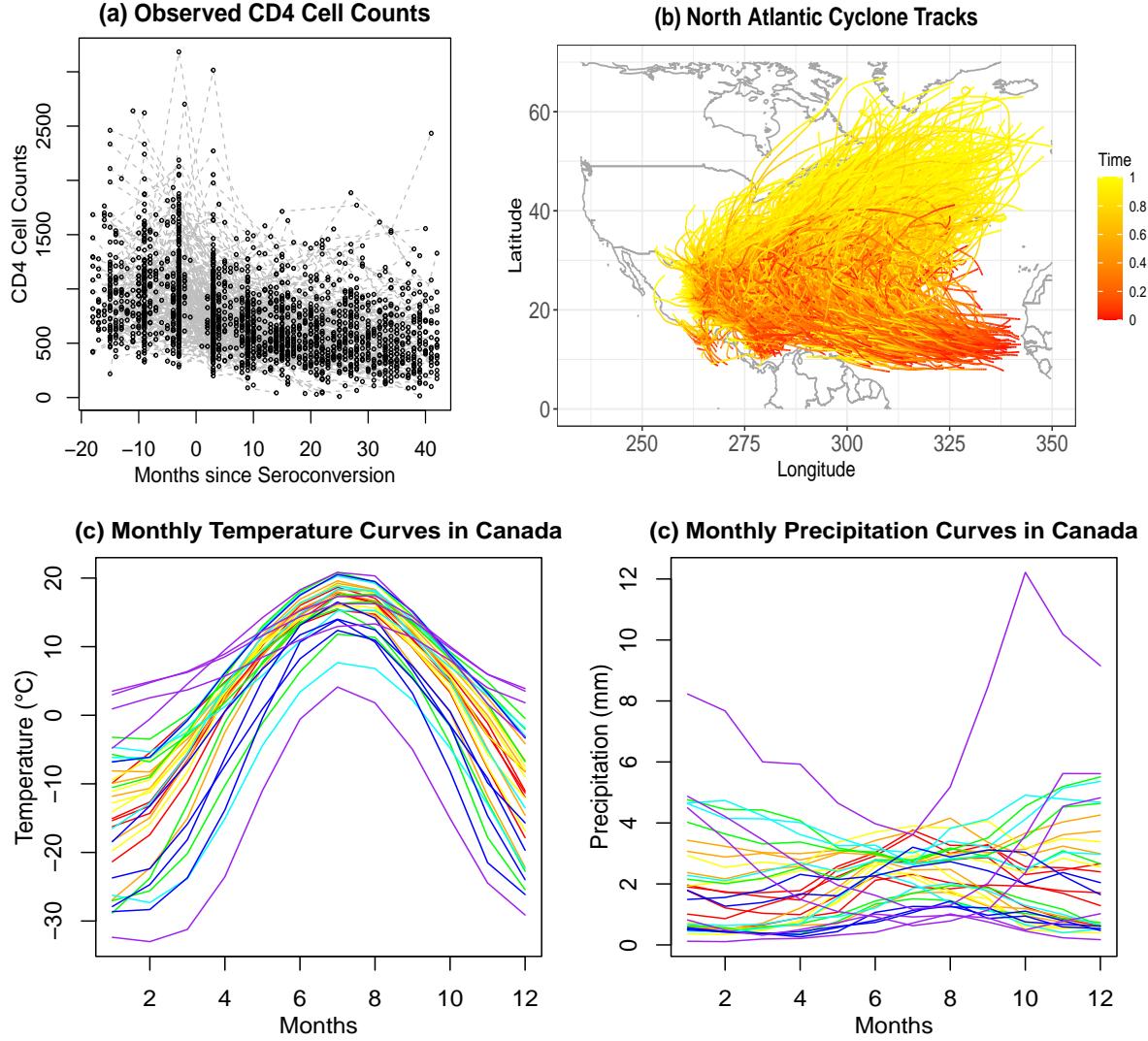


Figure 1: (a) shows the observed CD4 counts for 366 subjects during months -18 to 42 post seroconversion, (b) shows 1873 North Atlantic cyclone tracks recorded from 1851 to 2021, and (c) shows the rainbow plots of monthly temperature and precipitation curves at 35 different locations in Canada averaged over the period from 1960 to 1994. The orderings are based on the modified simplicial band depth ([López-Pintado et al. 2014](#)) of Canada temperature.

alization tools as well as for the exploratory analysis procedures applied to the data. During the past two decades, much effort has been made to find effective methods for exploratory functional data analysis (EFDA) and for estimating mean and covariance functions ([Yao et al. 2005, Wang et al. 2016, Happ & Greven 2018](#)). A series of methods and tools have

been developed, along with the proliferation of statistical models and inference techniques for functional data (Ramsay & Silverman 2005, Horváth & Kokoszka 2012, Wang et al. 2016).

This paper reviews the novel data mining methodology and visualization tools specifically used for FDA as an initial step prior to diving into modeling and statistical inference analysis. As compared to a case study of the geometric features of the internal carotid artery (2009), we introduce exploratory analysis with novel visualization tools and methods of clustering and classification. In addition, we use the univariate sparse CD4 data, the bivariate irregular North Atlantic cyclone track data, and the bivariate dense Canadian weather data to illustrate different methods. In contrast to the review of Wang et al. (2016), which concerns the general analysis of univariate functional data, we cover the EDA of p -dimensional ($p \in \mathbb{Z}_+$) functional data, where the measurement index per subject can vary. Hence, univariate functional data correspond to the special case of $p = 1$, and samples with dense time grids correspond to the special case of identical measurement indexes per subject.

The rest of the paper is organized as follows. Section 2 summarizes descriptive statistics for functional data. Section 3 proposes current tools for visualizing the observed functional data intuitively. Section 4 displays visualization tools featuring the descriptive statistics of functional samples. Section 5 presents several methods for functional data clustering and classification of dense and sparse functional data, separately. Section 6 concludes the paper with a summary and discussion.

2 Notations and Descriptive Statistics

In this section, we focus on the mathematical definitions and basic descriptive statistics of functional data, with an emphasis on the case in which the dataset is contaminated by

abnormal subjects.

2.1 Notation

A functional random variable \mathbf{Y} (Hsing & Eubank 2015) is a random vector with values in an infinite-dimensional space. Specifically, we view the p -variate ($p \in \mathbb{Z}_+$) functional data as paths of a p -variate stochastic process, taking values in some Hilbert space \mathcal{H} , such as the space of square-integrable functions defined on some bounded and closed interval \mathcal{T} . That is, $\mathcal{H} := L^2(\mathcal{T}) \times \cdots \times L^2(\mathcal{T})$. When $p = 1$, we return to univariate functional data.

Without loss of generality, we allow each marginal random vector in a p -variate stochastic process $\mathbf{Y}(\mathbf{t})$ to be defined at different indexes, that is, $\mathbf{Y}(\mathbf{t}) = (Y^{(1)}(t^{(1)}), \dots, Y^{(p)}(t^{(p)}))^\top$ with $\mathbf{t}^\top := (t^{(1)}, \dots, t^{(p)}) \in \mathcal{T} := \mathcal{T}_1 \times \cdots \times \mathcal{T}_p$. Note that \mathbf{t} is a p -dimensional vector, with its element $t^{(j)}$ being a random time and independent of all other random variables. Each element $Y^{(j)}(t^{(j)})$ ($j = 1, \dots, p$) is defined on the domain \mathcal{T}_j , where the \mathcal{T}_j s are compact sets in \mathbb{R} with finite Lebesgue measure. Briefly speaking, $Y^{(j)}(t^{(j)}) : \mathcal{T}_j \rightarrow \mathbb{R}$ is assumed to be square-integrable in \mathcal{T}_j , expressed as $L^2(\mathcal{T}_j)$. Then, we consider the p -dimensional functional data $\mathbf{Y} = \{\mathbf{Y}(\mathbf{t})\}_{\mathbf{t} \in \mathcal{T}}$ as sample paths of stochastic process $\mathbf{Y}(t)$, and we have $\mathbf{Y} \in \mathcal{H}$, where the space $\mathcal{H} := L^2(\mathcal{T}_1) \times \cdots \times L^2(\mathcal{T}_p)$.

In the following, let $\mathbf{Y}_1, \dots, \mathbf{Y}_N$ be a set of independent observations of \mathbf{Y} . In practice, we observe the functions $\mathbf{Y}_i(\mathbf{t})$ ($i = 1, \dots, N$) with error $\boldsymbol{\epsilon}_i(\mathbf{t}) = (\epsilon_i^{(1)}(t^{(1)}), \dots, \epsilon_i^{(p)}(t^{(p)}))^\top$, and the element $\epsilon_i^{(j)}(t^{(j)})$ are i.i.d. random variables with zero means. Moreover, the functions $\mathbf{Y}_i(\mathbf{t}_i)$ are observed on irregular finite grids at the subject and element level, that is, the j th ($j = 1, \dots, p$) element $t_i^{(j)}$ of \mathbf{t}_i ($\mathbf{t}_i \in \mathcal{T}$) can vary for each curve. Let the observed functions with measurement errors and the sparseness be $\tilde{\mathbf{Y}}_i(\mathbf{t}_i)$, such that $\tilde{\mathbf{Y}}_i(\mathbf{t}_i) = \mathbf{Y}_i(\mathbf{t}_i) + \boldsymbol{\epsilon}_i(\mathbf{t}_i)$.

2.2 Moment-Based Methods

We will consider a collection of functional data, $\{\mathbf{Y}_1, \dots, \mathbf{Y}_N\}$, consisting of N independent subjects observed at finite time points. The sampling schedule could vary in both location and number for each subject. Based on $\mathbf{Y}_i(\mathbf{t}_i)$ ($i = 1, \dots, N$) and $\mathbf{t}_i = (t_i^{(1)}, \dots, t_i^{(p)})^\top \in \mathcal{T}$, we define $\boldsymbol{\mu}(\mathbf{t}) := \text{E}\{\mathbf{Y}(\mathbf{t})\}$ as the mean function $\boldsymbol{\mu}$ evaluated at time \mathbf{t} , with the element estimation $\hat{\mu}^{(j)}(t^{(j)}) = \text{E}\{Y^{(j)}(t^{(j)})\} = \frac{\sum_{i=1}^N Y_i^{(j)}(t_i^{(j)}) \mathbf{1}(t_i^{(j)} = t^{(j)})}{\sum_{i=1}^N \mathbf{1}(t_i^{(j)} = t^{(j)})}$ for $l = 1, \dots, L$ and $j = 1, \dots, p$.

When functional data all have common and finite grid points, the number of observations at each grid point is equal to the number of subjects. However, when functional data are observed on irregular grids, the number of observations at each grid point varies and is imbalanced. It may be practical to count the number of observations in each bin rather than at each grid point. To obtain the whole curve, one can simply apply smooth interpolation or nonparametric smoothing methods, e.g., kernel smoothing (Wand & Jones 1995), local polynomial smoothing (Fan & Gijbels 1996), or spline smoothing (Wang 2011).

For $\mathbf{s}, \mathbf{t} \in \mathcal{T}$, we define the matrix of covariances $\mathbf{C}(\mathbf{s}, \mathbf{t}) := \text{cov}\{\mathbf{Y}(\mathbf{s}), \mathbf{Y}(\mathbf{t})\}$ with elements $C_{ij}(s^{(i)}, t^{(j)}) := \text{cov}\{Y^{(i)}(s^{(i)}), Y^{(j)}(t^{(j)})\}$ for $s^{(i)} \in \mathcal{T}_i$ and $t^{(j)} \in \mathcal{T}_j$. Likewise, the pointwise covariance function can be estimated as

$$\hat{C}_{ij}(t_k^{(i)}, t_l^{(j)}) = \frac{\sum_{n=1}^N \{X_n(t_n^{(i)}) - \hat{\mu}^{(i)}(t_k^{(i)})\} \{X_n(t_n^{(j)}) - \hat{\mu}^{(j)}(t_l^{(j)})\} \mathbf{1}(t_n^{(i)} = t_k^{(i)}, t_n^{(j)} = t_l^{(j)})}{\sum_{n=1}^N \mathbf{1}(t_n^{(i)} = t_k^{(i)}, t_n^{(j)} = t_l^{(j)})},$$

and the whole surface of the covariance function can be obtained by smoothing the three-dimensional (3D) scatterplot. Common smoothing methods for sparse functional data include multivariate functional principal component analysis (MFPCA) (Happ & Greven 2018) and tensor-product splines (Cai & Yuan 2010, Xiao et al. 2018, Li et al. 2020). Practically, these records are often assumed to be contaminated by measurement errors, and we refer the

readers to the article by [Wang et al. \(2016\)](#) for a comprehensive review of the estimation of mean and covariance functions in such a scenario.

2.3 Robust Methods

Functional data can be contaminated by abnormal subjects, also known as outliers, in a similar manner to univariate or multivariate data. Outliers may severely bias the aforementioned moment-based estimators and, consequently, lead to incorrect inference results. Hence, it is desired to develop methods that could eliminate the influence of outliers and summarize functional data robustly.

For univariate data, order-statistics and ranks induced naturally by the order of scalars on the real line are commonly used to design robust analysis methods, whereas for functional data, such a natural ranking is not available. During the past two decades, the idea of data depth, initially proposed to sort multivariate data, has been generalized to functional data. Specifically, a functional depth, taking values in $[0, 1]$, maps functional data as scalars and assigns larger depth values to central ones and smaller depth values to the more outward ones. Consequently, these scalars provide a ranking criterion for functional data in the data cloud from the center outward.

Commonly implemented depth notions for dense univariate functional data include but are not limited to band depth (BD) and modified band depth (MBD) ([López-Pintado & Romo 2009](#)), half-region depth and modified half-region depth (HRD and MHRD) ([López-Pintado & Romo 2011](#)), extremal depth ([Narisetty & Nair 2016](#)), functional tangential angle pseudo-depth (FUNTA) ([Kuhnt & Rehage 2016](#)) and its robustified version, order extended integrated depth ([Nagy et al. 2017](#)), total variation depth (TVD) ([Huang & Sun 2019](#)), and elastic depths ([Harris et al. 2021](#)). For dense multivariate functional data, available depth

notions include combinations of univariate functional depth measures (Ieva & Paganoni 2013), simplicial band depth (SBD) and modified simplicial band depth (MSBD) (López-Pintado et al. 2014), multivariate functional halfspace depth (MFHD) (Claeskens et al. 2014), and multivariate FUNTA pseudo-depth and its robustified version (Kuhnt & Rehage 2016).

For sparse univariate functional data, López-Pintado & Wei (2011) first proposed a model-based consistent procedure for estimating the depths based on the estimated curves on regular grids. Then, Sguera & López-Pintado (2021) proposed a new depth that enables the curve estimation uncertainty to be incorporated into the depth analysis. Those two depth notions have been extended to sparse multivariate functional data by Qu & Genton (2022), who also compared their ranking performances with simulations. In a recent study, Qu et al. (2022) introduced a novel framework for multivariate functional depths, specifically designed for sparse multivariate functional data and eliminating the need for curve estimation. This new depth concept, termed “global depth”, distinguishes itself from previous approaches by handling sparse functional data directly. The authors have demonstrated how the procedures for MFHD and multivariate extremal depth (an extension of extremal depth) can be adapted to their global depth framework.

Functional depths provide a natural basis for defining the median, extremes, and quantiles of functional data. Fraiman & Muniz (2001) defined the functional median as the deepest observation, i.e., the sample with the largest depth value, denoted as $\mathbf{M} = \arg \max_{\mathbf{X}} D(\mathbf{X}, F_{\mathbf{X}})$, where $D(\mathbf{X}, F_{\mathbf{X}})$ denotes the depth of a random function \mathbf{X} with respect to its distribution $F_{\mathbf{X}}$.

The functional version of the α trimmed mean, $\boldsymbol{\mu}_{\alpha}$, is defined as the average of the deepest

$1 - \alpha$ proportion of subjects,

$$\boldsymbol{\mu}_\alpha = \frac{E(\mathbf{X} \mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}, F_{\mathbf{X}})))}{E(\mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}, F_{\mathbf{X}})))},$$

where $\mathbf{1}_A(x) = 1$ for $x \in A$ and zero otherwise, and $E(\mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}, F_{\mathbf{X}}))) = 1 - \alpha$. The empirical definitions of the two statistics can be expressed as

$$M_N = \arg \max_{i=1, \dots, N} D(\mathbf{X}_i, F_{\mathbf{X}, N}) \quad \text{and} \quad \hat{\boldsymbol{\mu}}_{\alpha, N} = \frac{\sum_{i=1}^N \mathbf{X}_i \mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}_i, F_{\mathbf{X}, N}))}{\sum_{i=1}^N \mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}_i, F_{\mathbf{X}, N}))}.$$

Similarly, the α trimmed covariance function can be defined as

$$C_\alpha(s, t) = \frac{E[\{\mathbf{X}(s) - \boldsymbol{\mu}_\alpha(s)\}\{\mathbf{X}(t) - \boldsymbol{\mu}_\alpha(t)\}\mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}, F_{\mathbf{X}}))]}{E(\mathbf{1}_{[\beta, \infty)}(D(\mathbf{X}, F_{\mathbf{X}})))},$$

and its empirical version can be derived by substituting the statistics with their respective estimators. Another concept related to the ranking of functional data is the central region (López-Pintado & Romo 2009, Sun & Genton 2011, Narisetty & Nair 2016, Myllymäki et al. 2017), which is defined as

$$C_{1-\alpha} = \{\mathbf{X} \in L_2(\mathcal{T}) : X_L^{(j)}(t) \leq X^{(j)}(t) \leq X_U^{(j)}(t), \forall t \in \mathcal{T}, j = 1, \dots, p\},$$

where \mathbf{X}_L and \mathbf{X}_U are lower and upper α -envelope functions, $\mathbf{X}_L = \inf\{\mathbf{X} \in L_2(\mathcal{T}) : D(\mathbf{X}, F_{\mathbf{X}}) > \alpha\}$ and $\mathbf{X}_U = \sup\{\mathbf{X} \in L_2(\mathcal{T}) : D(\mathbf{X}, F_{\mathbf{X}}) > \alpha\}$, respectively.

Similarly, the global envelope (Myllymäki & Mrkvička 2019) is a band that the functional data \mathbf{X} falls outside with the probability α . For the univariate functional data $\mathbf{X} = \{X(t), t \in \mathcal{T}\}$, we define $\mathbf{X}_{low} = \{X_{low}(t), t \in \mathcal{T}\}$ and $\mathbf{X}_{upp} = \{X_{upp}(t), t \in \mathcal{T}\}$ as the lower and upper bounds, respectively. Then, the global envelope $[X_{low}(t), X_{upp}(t)]$ is

defined as

$$P(X(t) \notin [X_{low}(t), X_{upp}(t)] \text{ for } t \in \mathcal{T}) = \alpha.$$

3 Simple Visualization

Visualization ([Friedman & Stuetzle 2002](#)) has long been a component of great importance to EDA, and many visualization tools are widely used as routine steps in the analysis procedure. For instance, the histogram of a univariate dataset shows a rough sense of the density of its underlying distribution, the scatter plot of a bivariate dataset presents the locations of the data points on a two-dimensional plane to provide some idea of the relation between the two variables, and a heatmap shows the magnitude of an object as color in two dimensions. Similar demands in FDA motivate researchers to develop new graphical tools.

Here, we highlight several reasonably simple tools that have proved useful in the literature ([Hyndman & Shang 2010](#), [Hubert et al. 2015](#), [Wrobel & Goldsmith 2016](#)). We will consider the Canadian weather dataset from Figure 1 (c) as one instance. This dataset includes monthly recorded temperature and precipitation curves for 35 stations in Canada over the period from 1960 to 1994.

3.1 Spaghetti Plot and Rainbow Plot

A spaghetti plot ([Allen 2019](#)) is a simple visualization that assigns a distinct color to each subject, making it easy to track movement for data with small samples. However, such a plot may look messy when used to visualize big functional data. The rainbow plot, proposed by [Hyndman & Shang \(2010\)](#), can be regarded as an improvement of the spaghetti plot. As a visualization of all the curves, it adds a data ordering feature and colors the samples based

on the ordering, using the rainbow palette. The order can reflect time, data depth, data density, or another index.

In Figure 1 (c), the Canadian temperature curves are ordered with the MSBD from the median to the extreme, and the curves are labeled from red to purple in the rainbow palette. We can see that the red group represents the median tendency of the temperature and precipitation over the course of a year, whereas the purple group includes data from some stations with high temperatures and high precipitation during winter and some stations with low temperatures and low precipitation all year round.

3.2 Heatmap

A heatmap ([Hubert et al. 2015](#)) represents different values by using a system of color-coding. In FDA, a $n \times m$ heatmap is suitable for showing a functional dataset consisting of n subjects recorded on m common design points. We visualize the data with a heatmap in Figure 2. For instance, each cell in Figure 2 (a) represents the estimated temperature at one station in a specific month, each row represents the monthly temperature curve for a station, and each column represents the average temperature at 35 stations in a particular month. Some abnormal information can be easily detected through the heatmap. Figure 2 (a) shows that Victoria and Vancouver have persistent high temperatures between April and October, whereas temperatures in Resolute and Iqaluit remain below 10 degrees Celsius almost all year round. Pr. Rupert has monthly precipitation up to 6 mm, except between April and August, whereas the other stations have monthly precipitation of less than 6 mm.

3.3 Interactive Plots

Several packages have been developed to generate interactive visualization for functional data. An interactive plot retains the advantages of both visual and numerical illustration of data, i.e., it is intuitive as well as accurate. The interaction can be achieved in many ways, e.g., by showing the associated records at the locations indicated by the cursor, zooming in or out, or interacting between different plots. [Wrobel et al. \(2016\)](#) proposed using the *refund.shiny* package that creates interactive graphics for FDA. The *refund.shiny* package ([Chang et al. 2015](#)) relies on the *shiny* package to generate such an interactive user interface. Another commonly used tool is the *plotly* package ([Sievert et al. 2018](#)), which produces interactive plots with two or three dimensions in combination with a web portal.

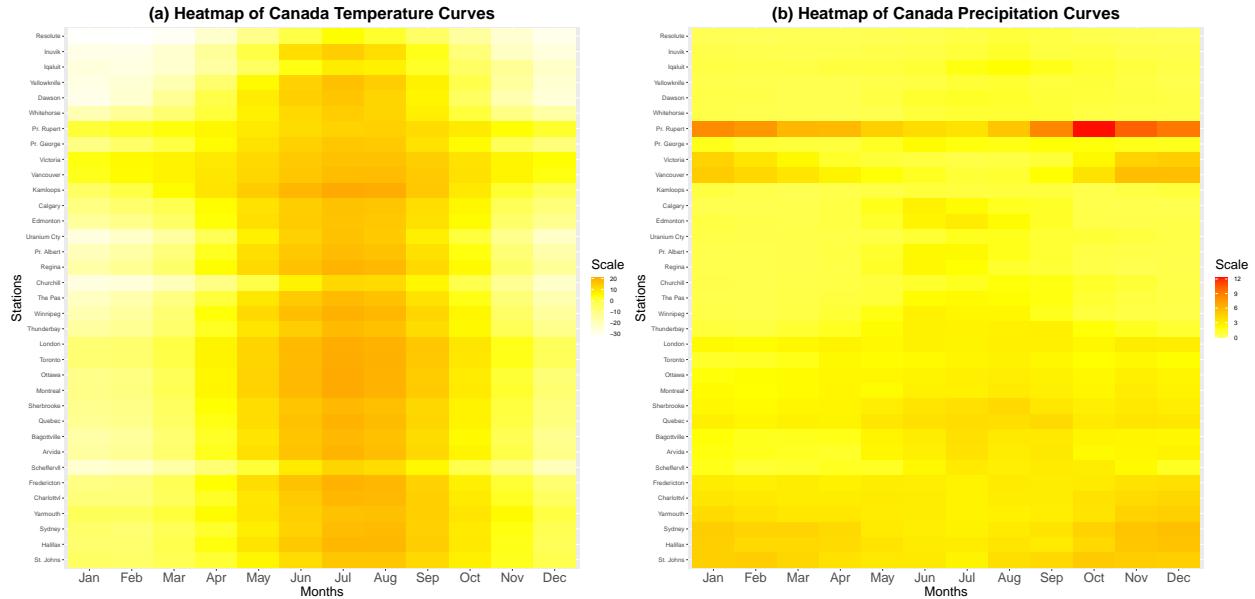


Figure 2: (a) is the heatmap of 35 Canada monthly average temperature curves over the period from 1960 to 1994, and (b) is the heatmap of 35 Canada monthly average precipitation curves over the period from 1960 to 1994.

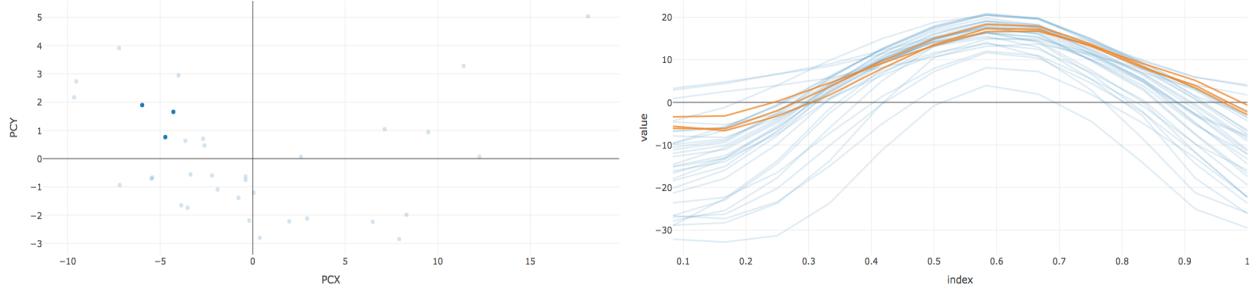


Figure 3: An illustration of an interactive functional principal component (FPC) plot generated by the *refund.shiny* package. The left panel shows the observed score scatterplot for selected FPCs of Canadian temperature, and the right panel shows fitted values for subjects, where subjects selected in the left panel are in orange and the other subjects are in blue.

3.4 Animations

Animation or video is another powerful tool for enhancing still figures that can visualize the dynamic evolution of data. Genton et al. (2015) proposed the term *visuanimation* for referring to visualization through animations, and they explored the utility of animation in various perspectives of statistics. Castruccio et al. (2019) illustrated the predicted global temperature data, which can be regarded as functional data, varying spatially and temporally via a 3D virtual-reality movie, and they developed a mobile application that enables users to watch the movie interactively.

4 Visualization with Structure Information

Many visualization tools for classical data have been developed that feature descriptive statistics. For instance, the boxplot (Tukey 1977) of a univariate dataset illustrates the structure of the dataset by showing its descriptive statistics, e.g., the median, quartiles, extreme values, and possible outliers; whereas the bagplot (Rousseeuw et al. 1999) of a bivariate dataset presents the deepest data, the deepest 50% of the data points, and possible

outliers under the ranks given by the halfspace depth (Tukey 1975). In FDA, functional data can be transformed to notions of depth or outlyingness (Dai & Genton 2019) from the center outwards (see Subsection 2.3) for visualization and outlier detection. Hence, we will introduce visualization tools that contain structure information of raw data.

4.1 Visualization Based on Ranking Information

Hyndman & Shang (2010) first proposed several visualization tools for smoothed functional data, such as the functional bagplot and functional highest density region (HDR) boxplot available in the R package *rainbow* (Shang & Hyndman 2019). The functional bagplot is based on the bivariate bagplot of Rousseeuw et al. (1999). It first applies the first two robust principal component scores to the bivariate bagplot as an auxiliary tool to rank the observations and detect outliers. Then, it displays the median curve, the 50% inner region, and the 99% fence. Curves that are partially outside these regions are identified as outliers. The functional HDR boxplot is a mapping of the bivariate HDR boxplot (Hyndman 1996) of the first two robust principal component scores to the functional curves. In contrast to the functional bagplot, this method displays curves with high HDRs. Specifically, it focuses on curves whose first two principal component scores correspond to the 50% inner and 99% outer bivariate HDRs. Additionally, it identifies outliers as points that are excluded from the 99% outer HDR. The visualization tools mentioned above use either time-ordering or the first two robust principal component scores to arrange the curves.

The functional boxplot, as proposed by Sun & Genton (2011), is a data visualization technique used to summarize the distribution and features of a set of functional data. It uses the functional depth and highlights the central quantiles and possible outliers. Analogous to the classical boxplot, there are four descriptive statistics in the functional boxplot (see Figure

[4](#)): the envelope of the 50% central region, the median curve, the outliers, and the maximum non-outlying envelope. An observation is flagged as an outlier if its measurement at any grid point is outside a constant factor times the range at the central region. The constant factor is set to be 1.5 under the assumption that observations at each index are independent and identically distributed and that they follow a normal distribution. The functional boxplot is generalized to other types of boxplots to suit functional data with additional characteristics. We can categorize the various functional boxplots as follows: those with more descriptive statistics, those dealing with spatio-temporal data, those with missing data, and those with more general data objects.

In the first category, the enhanced functional boxplot ([Sun & Genton 2011](#)), the double-fence functional boxplot ([Serfling & Wijesuriya 2017](#)), and the two-stage functional boxplot ([Dai & Genton 2018a](#)) were proposed to underline more features. For instance, the enhanced functional boxplot provided 25% and 75% central regions on the basis of the functional boxplot and the two-stage functional boxplot; the double-fence functional boxplot included an additional fence of 0.5 interquartile regions, enhancing its ability to identify specific shape and location outliers; and the two-stage functional boxplot implemented the directional outlyingness ([Dai & Genton 2019](#)) first, colored the detected outliers in green, then applied the remaining curves to the procedures in the functional boxplot.

In the second category, the adjusted functional boxplot ([Sun & Genton 2012](#)) and surface boxplot ([Genton et al. 2014](#)) were proposed for use with spatio-temporal data. The spatio-temporal data can be viewed as a temporal curve at each spatial location or as a spatial surface at each time. In the former case, correlations need to be considered across locations. Hence, [Sun & Genton \(2012\)](#) proposed the adjusted functional boxplot, which flexibly selects the constant factor to control the probability of correctly detecting no outliers. In the

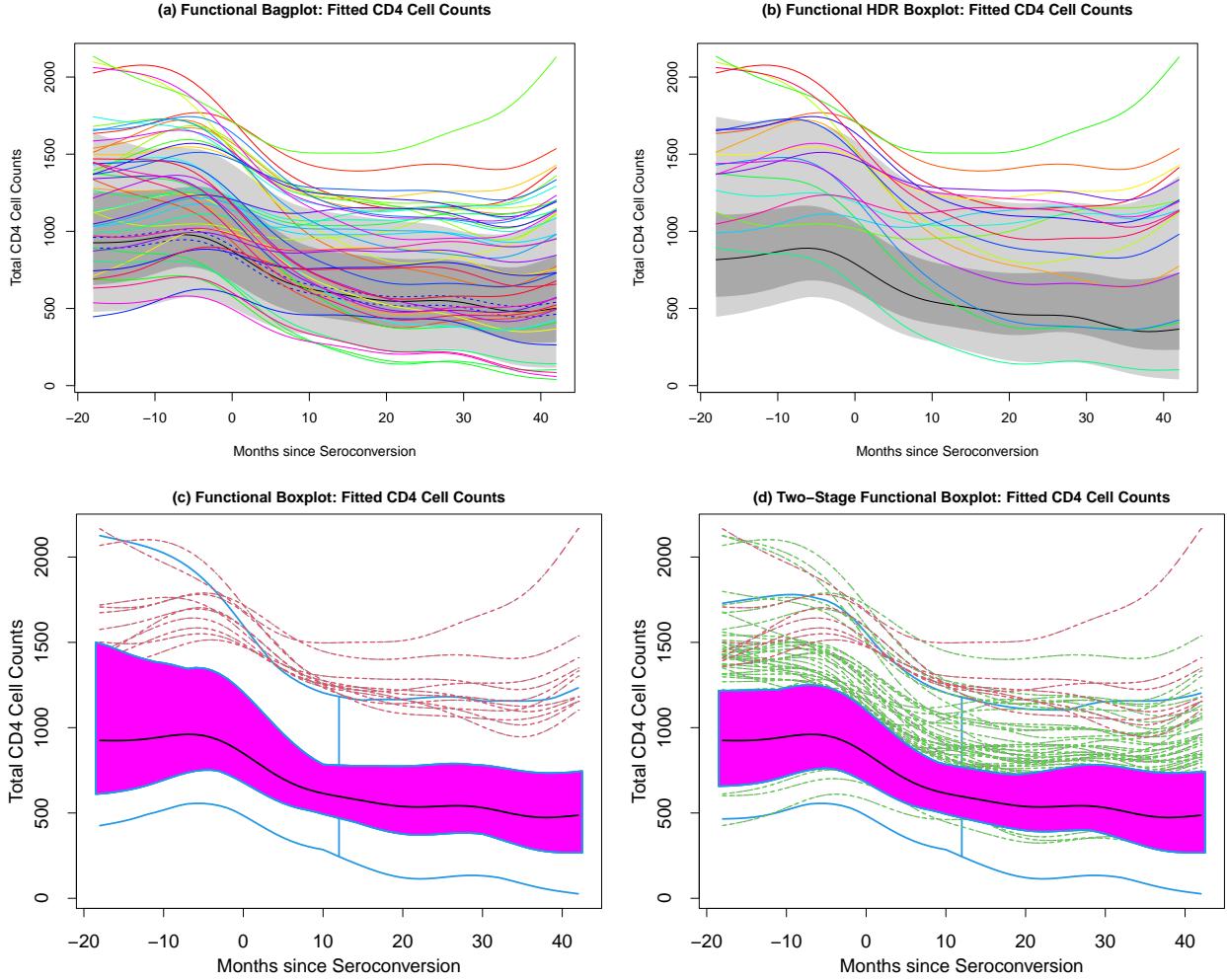


Figure 4: Comparisons of the functional bagplot, the functional HDR boxplot, and the functional boxplots of the fitted CD4 cell counts from bootstrap MFPCA (Qu & Genton 2022): (a) the functional bagplot, (b) the functional high-density region (HDR) boxplot, (c) the functional boxplot, and (d) the two-stage functional boxplot.

aforementioned work, Genton et al. (2014) extended the concept of MBD to modified volume depth specifically for image data. This extension enabled them to introduce a surface boxplot, which facilitates the visualization of image subjects based on the modified volume depth. Similarly, the same four descriptive statistics can also be established by using the modified volume depth.

In the third category, the sparse functional boxplot and the intensity sparse functional

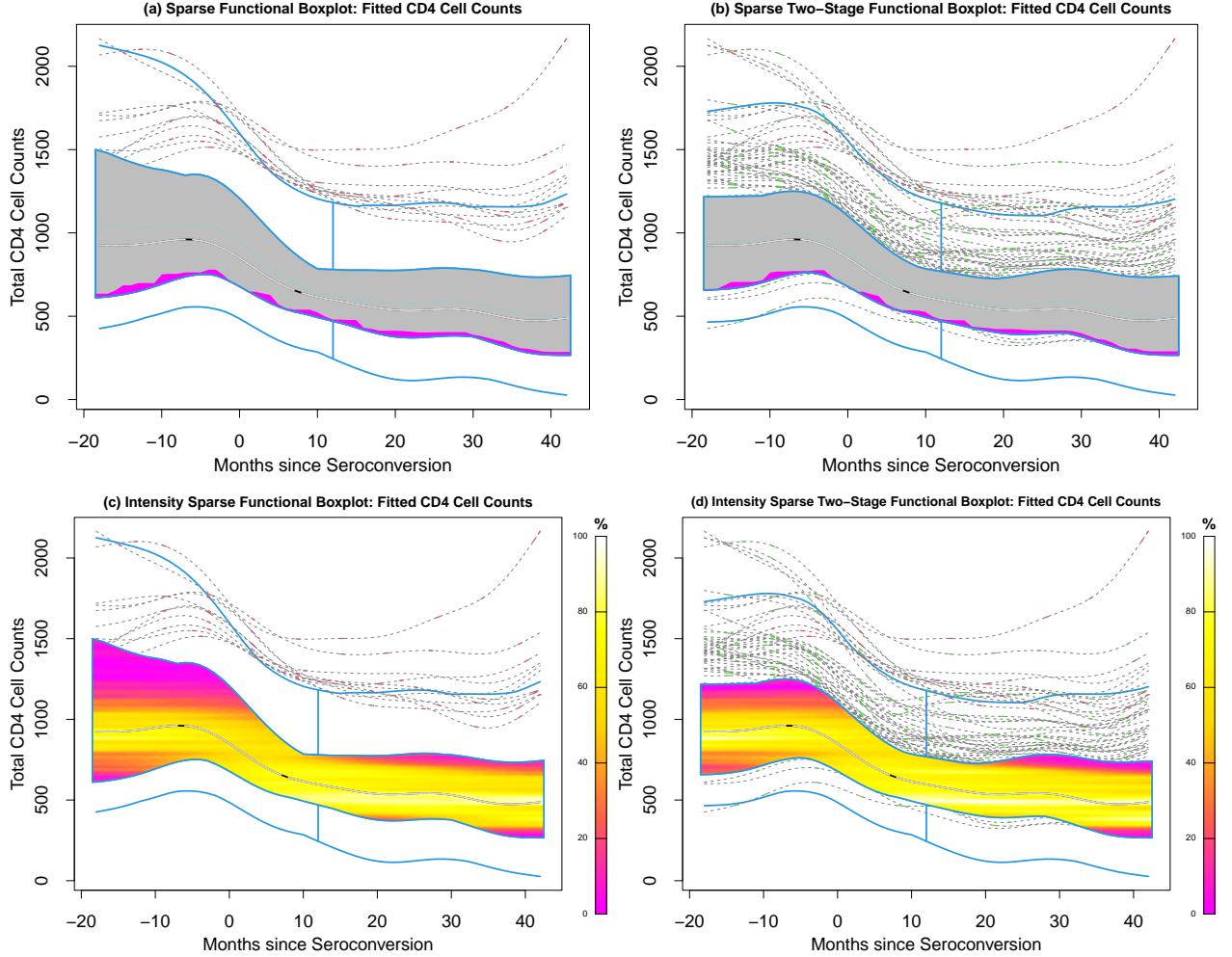


Figure 5: Functional boxplot and its variations when missing values exist, taking the instance of the fitted CD4 cell counts from bootstrap MFPCA (Qu & Genton 2022). The left column includes (a) the sparse functional boxplot and (c) the intensity sparse functional boxplot. The right column includes (b) the sparse two-stage functional boxplot and (d) the intensity sparse two-stage functional boxplot.

boxplot (Qu & Genton 2022) were proposed for visualization. Data reconstruction is required with MFPCA (Happ & Greven 2018). In addition to the descriptive statistics in the functional boxplot, the sparse functional boxplot displays the smooth sparseness proportion within the 50% central region, and the intensity sparse functional boxplot displays the intensity of the smooth sparseness within the 50% central region. Usually, the directional

outlyingness (Dai & Genton 2019) and sparse functional boxplots are combined to form the sparse two-stage functional boxplot and the intensity sparse two-stage functional boxplot for visualization and outlier detection (see Figure 5). Furthermore, sparse functional boxplots have been extended to the simplified sparse functional boxplot (Qu et al. (2022)) and the simplified intensity sparse functional boxplot without data reconstruction. The simplified visualization tools are based on the global multivariate functional depths (Qu et al. 2022), which are applied to the sparse multivariate functional directly without data reconstruction.

The fourth category includes other natural extensions of the functional boxplot for data expressed as sets, curves, paths, or trajectories. Whitaker et al. (2013) defined the set BD and introduced a contour boxplot for visualization and exploration of ensembles of contours or level sets of functions. Mirzargar et al. (2014) generalized the BD for curves and proposed the curve boxplot. Hong et al. (2014) introduced a weighted functional boxplot for use when data objects become shapes and images. Raj et al. (2017) proposed the graph-simplex band depth and developed a visualization tool path boxplot. Yao et al. (2020) developed a trajectory boxplot (Figure 6) for visualizing and exploratory analysis of trajectories that show variation in longitude and latitude through time.

The functional boxplot has some shortcomings, such as the loss of functional interpretation in the envelopes of the 50% central region and the non-outlying region and its inapplicability to functional observations under hidden temporal warping variability. Therefore, Xie et al. (2017) decomposed observation variation in functional data into three main components: amplitude, phase, and vertical translation based on the curve registration (Srivastava et al. 2011). They constructed a different visualization for each element based on the median, two quartiles, and extreme observations. They also proposed identifying outliers based on those three components and visualizing the amplitude or phase outliers through

the phase-versus-amplitude distance plot.

4.2 Centrality Decomposition Plots

Another set of visualizations, specifically for outlier detection, are usually based on ranking criteria such as statistical depth or outlyingness. The outlier detection visualization tools that we introduce here can be separated into those specifically for univariate functional data and those for multivariate functional data (univariate functional data are usually special cases).

For univariate functional data, [Arribas-Gil & Romo \(2014\)](#) proposed using outliergrams to visualize and detect shape outliers in functional data by exploiting the relation between MBD and the modified epigraph index (MEI, [López-Pintado & Romo 2011](#)). Through a novel decomposition of the total variation depth, proposed by [Huang & Sun \(2019\)](#), we can easily detect shape outliers via the boxplot of the modified shape similarity (MSS). Thereafter, the magnitude outliers can be seen among the remaining observations with the functional boxplots.

For multivariate functional data, [Hubert et al. \(2015\)](#) discussed amplitude and shape outliers and proposed various functional outlier maps based on the notion of outlyingness and depth in the multivariate functional data, e.g., adjusted outlyingness (AO) and skew-adjusted projection depth (SPD). Furthermore, they exploited the relation between AO and SPD. They constructed the centrality-stability plot in which the amplitude outliers lie in the upper-right region and the shape outliers lie in the right region. [Rousseeuw et al. \(2018\)](#) proposed a robust notion of outlyingness, directional outlyingness (DO), and this can be applied in the univariate or multivariate setting. Based on the DO in the univariate setting, they defined the average of outlyingness as the functional directional outlyingness (FO)

and measured the variability of its DO (VO). Then, they developed a graphical tool called the functional outlier map (FOM), which is a scatterplot of (FO, VO). Shift outliers, local outliers, and global outliers can be detected and displayed in different domains in FO. Based on the relation between mean directional outlyingness (MO) and VO, [Dai & Genton \(2018b\)](#) proposed a new graphical tool, the magnitude-shape (MS) plot, to illustrate the centrality of curves comprehensively. They also generalized the outliergram to the bivariate outliergram for outlier detection in bivariate functional data, according to a quadratic relation between FO and MO. However, the bivariate outliergram is limited to bivariate functional data and is less efficient than the MS plot at measuring the centrality of curves. [Yao et al. \(2020\)](#) introduced wiggleness of directional outlyingness (WO) to detect outliers, and constructed the WO-MSBD plot which can distinguish shape outliers and magnitude outliers. The depth boxplot, introduced by [Harris et al. \(2021\)](#), is constructed on the elastic depths directly and serves as a half-boxplot. Its purpose is to identify potential amplitude and phase outliers. [Ojo et al. \(2023\)](#) proposed the magnitude-shape-amplitude plot (MSA-plot) based on fast massive unsupervised outlier detection (FastMUOD, [Ojo et al. 2022](#)).

5 Clustering and Classification

In the terminology of machine learning, functional data clustering is an unsupervised learning process, whereas functional data classification is a supervised learning procedure. Cluster analysis aims to group a set of data such that data objects are more similar within clusters than across clusters with respect to a given metric. In contrast, classification assigns a new data object to a predetermined group by using a discriminant function or classifier. The assumption of functional clustering is the existence of inhomogeneous groups, and the

goal is to estimate the membership of each group based on unlabeled observed data. Functional classification typically involves training data containing a functional predictor with an associated multiclass label for each data object. The discrimination procedure of functional classification is closely related to functional cluster analysis, even though the goals are different.

When the structures or centers of clusters can be established in functional data clustering, the criteria used for identifying clusters can also be used for classification. Clustering and classification are both valuable tools for EFDA. Detecting meaningful and homogeneous clusters can improve the estimation of some descriptive statistics, such as the mean or median. Additionally, correctly assigning unlabeled functional data to predetermined classes enables better understanding and prediction of updated observations.

5.1 Functional data clustering

The range of applications for functional data clustering is vast. For example, [Abramowicz et al. \(2017\)](#) applied a functional clustering method to study sediment data and to infer past environmental and climate changes. [Athanasiadis & Mrkvicka \(2019\)](#) analyzed financial time series by using functional clustering methods to identify different insurance penetration (IP) rate profiles in European markets. In general, the resulting clusters show high potential for data visualization and interpretation.

When dealing with functional data, similarities might take into account the characteristics of the curves, such as their shapes, magnitudes, or derivatives ([Hitchcock & Greenwood, 2015](#)). Broadly, we can classify the existing functional clustering methods as follows ([Jacques & Preda, 2014](#)): 1) raw data methods, 2) filtering methods, 3) adaptive methods, and 4) distance-based methods. Raw data methods represent a naive approach and might result

in high-dimensional vectorial clustering ([Bouveyron & Brunet-Saumard, 2014](#)). Filtering methods and adaptive methods use the basis expansion approach for functional data with a common basis for all of the data or a common basis per group, respectively. Some examples of clustering methods that use B-splines, Fourier basis, or functional principal component analysis have been described in detail by [Abraham et al. \(2003\)](#), [Serban & Wasserman \(2005\)](#), and [Shang \(2014\)](#). Lastly, distance-based methods quantify the similarity between clusters by computing distances for functional objects. Here, we focus on distance-based methods. For a useful review of filtering and adaptive methods, we refer the reader to the articles by [Jacques & Preda \(2014\)](#) and [Wang et al. \(2016\)](#).

There are two main objects in a distance-based clustering method: the similarity measure and the clustering algorithm. We need to define a *similarity or dissimilarity measure* between curves that will be highly related to the interpretation of the clusters. Usually, these measures are defined between two curves, $\{\mathbf{X}_i, \mathbf{X}_j\}$, where $\mathbf{X}_i = (X_i^{(1)}, \dots, X_i^{(p)})^\top$ and $\mathbf{X}_j = (X_j^{(1)}, \dots, X_j^{(p)})^\top$ are p -variate functional data, and we need a *clustering algorithm* to compute similarities across clusters, $C_1 = \{\mathbf{X}_1^1, \dots, \mathbf{X}_{n_1}^1\}$ and $C_2 = \{\mathbf{X}_1^2, \dots, \mathbf{X}_{n_2}^2\}$. Usually, the similarity measure can be defined using a distance between functions, $d(\mathbf{X}_i, \mathbf{X}_j)$. Natural choices for the distance are the L_1 , L_2 , or L_∞ distances, where

$$d_l(\mathbf{X}_i, \mathbf{X}_j) = \left(\frac{1}{p} \sum_{k=1}^p \int_{\mathcal{T}_k} (|X_i^{(k)}(t) - X_j^{(k)}(t)|)^l dt \right)^{1/l}$$

for $l = 1, 2, \infty$. If we consider the L_1 , L_2 , or L_∞ distance, then the resulting clusters are built of functions with similar shapes and magnitudes. If there is no interest in the similarity of magnitude, then the functions can be normalized and the total variation (TV) distance

used ([Alvarez-Esteban et al., 2016](#)):

$$d_{TV}(\mathbf{X}_i, \mathbf{X}_j) = 1 - \frac{1}{p} \sum_{k=1}^p \int_{\mathcal{T}_k} \min\{X_i^{(k)}(t), X_j^{(k)}(t)\} dt = \frac{1}{2p} \sum_{k=1}^p \int_{\mathcal{T}_k} |X_i^{(k)} - X_j^{(k)}| dt.$$

These distances might be more complex if we also include information about the derivative curve and define a similarity measure as a weighted combination of the distances i.e., $a_1 d(\mathbf{X}_i, \mathbf{X}_j) + a_2 d(\mathbf{X}'_i, \mathbf{X}'_j)$. Here, we assume that the curves \mathbf{X}'_i s are all independent. However, if the user is interested in clustering dependent curves, then a similarity measure can be proposed that uses the Spearman correlation or the rank correlation between functions ([Heckman & Zamar, 2000](#)). If these curves are linked to a time series trajectory, then a coherence-based distance might be useful too ([Euán et al., 2019](#)). In this setting, the correlation of the resulting clusters is high within each group but low across clusters. [Dai et al. \(2021\)](#) induced the dissimilarity matrix from functional ordering. The idea is to construct the set of functional differences, apply any functional depth (or ranking) notions to the above set, and define the similarity as one minus the depth.

However, those methods assume that functions are observed at a fixed set of points, and no sparseness exists. Elastic time distance was proposed by [Qu et al. \(2023\)](#) to address this issue. It is applicable to (multivariate) functional data with either identical or different time measurements per subject. The core idea is to build standard grid points and to interpolate measurements at standard grid points with the available observations. Assume curves \mathbf{X}_i and \mathbf{X}_j are p -variate multivariate functional data and that $\widetilde{\mathbf{X}}_i$ and $\widetilde{\mathbf{X}}_j$ are their interpolated observations based on procedures described by [Qu et al. \(2023\)](#), then

$$d_{ETD}(\mathbf{X}_i, \mathbf{X}_j) = \max_{k=1, \dots, T} \sqrt{\sum_{l=1}^p \{\widetilde{X}_i^{(l)}(t_k) - \widetilde{X}_j^{(l)}(t_k)\}^2}, t_k \in \mathcal{T}.$$

Once the similarity measure is chosen, we use a clustering algorithm that selects the groups of functions that are more similar in an “optimal” manner, i.e., members within each group are highly similar, but members across groups are highly dissimilar. The algorithms most commonly used for this purpose are the k -means and hierarchical clustering algorithms. [Ferraty & Vieu \(2006\)](#) introduced examples of hierarchical clustering using the L_2 distance between the functions and their second derivatives. [Ieva et al. \(2013\)](#) applied a k -means to identify clusters of electrocardiograph traces with a weighted distance between the curves and their first derivatives. Recently, [Euán et al. \(2018\)](#) proposed the hierarchical merger clustering algorithm. Its main contribution to classical hierarchical algorithms is the use of a representative member for each cluster. [Euán et al. \(2018\)](#) proposed using the TV distance in a hierarchical merger algorithm to cluster spectral density functions from ocean wave time series. [Euán & Sun \(2019\)](#) extended this method for general 2D directional spectra functions. Moreover, [Qu et al. \(2023\)](#) combined the elastic time distance and the original robust two-layer partition (RTLP) clustering algorithm to cluster multivariate functional curves. They also compared RTLP clustering with other algorithms, including the distance-based methods DBSCAN, k -means, and k -median and the model-based funHDDC algorithm ([Schmutz et al. 2020](#)).

Some real data applications might need a more robust clustering algorithm, especially if the data have a high noise level. Although some of the methods described previously in this section might separate possible outliers as single clusters, this is not true for all methods. In the presence of potential outliers, [Cuesta-Albertos & Fraiman \(2007\)](#) proposed a trimmed k -means clustering that results in a robust cluster procedure for functional data. Also, [Rivera-García et al. \(2019\)](#) applied the trimming technique to introduce a robust model-based clustering method for functional data. When data are misaligned, applying clustering

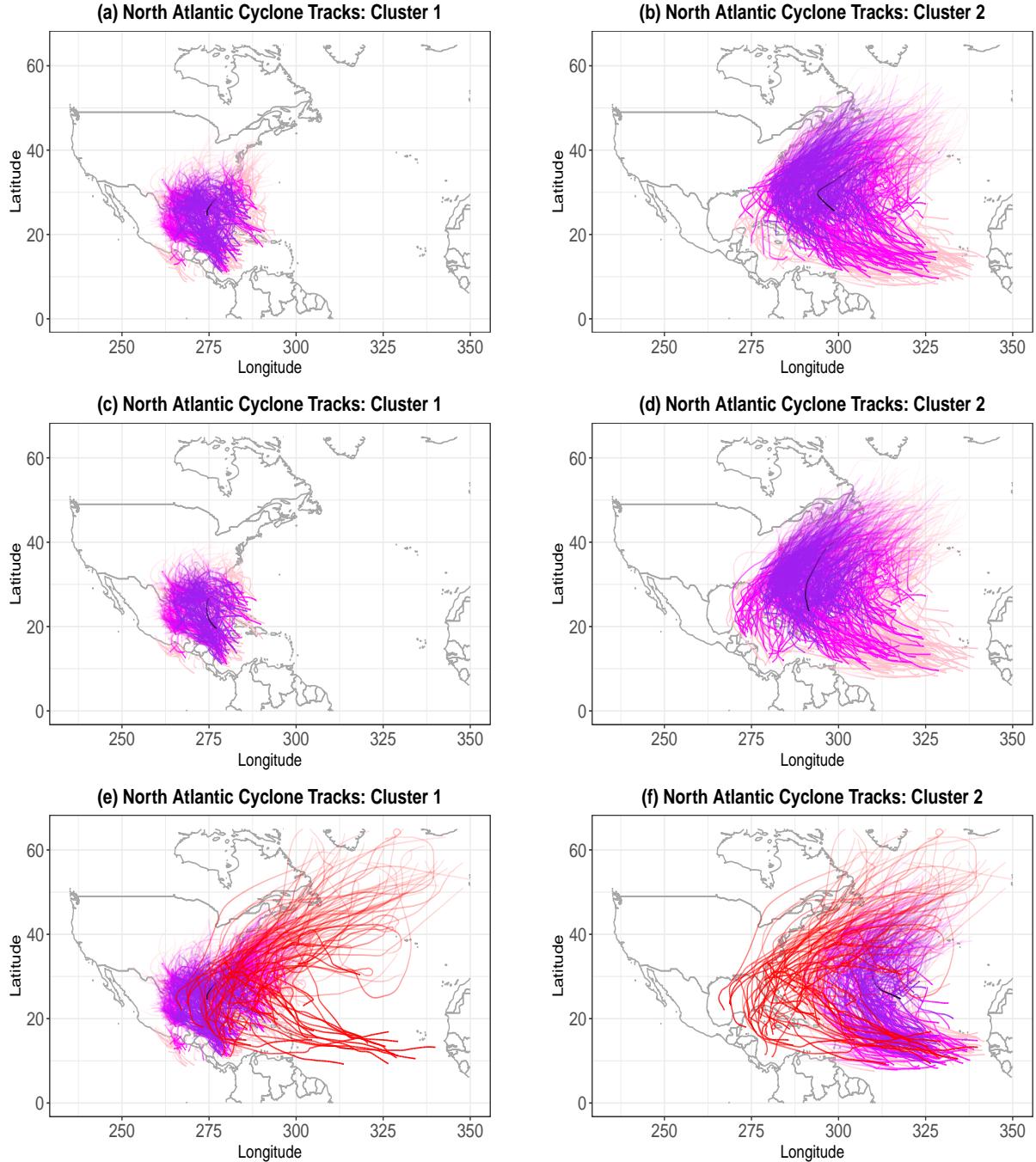


Figure 6: Cyclone trajectories obtained using our version of the trajectory boxplot (Yao et al. 2020). (a) and (b) are from k-medoids clustering, (c) and (d) are from hierarchical clustering, and (e) and (f) are from robust two-layer partition clustering. Black and red represent the median and outliers, respectively, and purple, magenta, and pink indicate the first, second, and third quartile curves, respectively.

methods directly might result in nonreasonable clustering structures. Sangalli et al. (2010) proposed an algorithm that considers the case in which curves are misaligned. De Micheaux et al. (2021), based on the curve depth, employed the original clustering algorithm (Jörnsten 2004) with slight modifications for unparameterized curves. The robust two-layer partition clustering, introduced by Qu et al. (2023), uses both a two-layer partition algorithm and a modified silhouette index. This approach is effective at distinguishing clusters and identifying potential outliers in terms of their magnitude and shape.

In general, a good strategy is to select the clustering method based on the research goal. We will illustrate this by using the bivariate hurricane trajectory data for the North Atlantic (see Figure 1 (b)). Because the trajectory data have various observations per subject, we apply the elastic time distance mentioned in Qu et al. (2023) and consider the following different clustering methods based on the interpolated data: 1) K -medoids clustering (Park & Jun 2009) (see Figure 6 (a)-(b)); 2) hierarchical clustering with the average as the linkage function (see Figure 6 (c)-(d)); and 3) robust two-layer partition clustering (see Figure 6 (e)-(f)). Each clustering method generates two clusters, but outliers are introduced only by the robust two-layer partition clustering.

5.2 Functional data classification

Whereas clustering seeks to find homogeneous clusters without knowledge of the true clusters, functional classification assigns a group membership to a new data object with a discriminant function or a classifier. To construct such a classifier, we assume that the observed data are $\{(\mathbf{X}_1, G_1), \dots, (\mathbf{X}_n, G_n)\}$, where $G_i \in \{1, 2, \dots, k\}$ is a categorical variable that indicates the class membership of curve \mathbf{X}_i . Then, for a new observed curve \mathbf{X}_0 , the goal of functional data classification is to assign the unknown class membership G .

Broadly, we can consider three different approaches to build the classifier (Wang et al., 2016): 1) regression-based functional classification, 2) functional discriminant analysis, and 3) depth-based classification. Regression-based classifiers assume that the response of the regression model is categorical, so that the assigned class membership corresponds to the value with the maximal probability. For instance, using a baseline functional logistic regression,

$$\log \frac{\mathbb{P}(G_0 = g | \mathbf{X}_0)}{1 - \sum_{j=1}^{k-1} \mathbb{P}(G_i = j | \mathbf{X}_i)} = \gamma_{0,g} + \int_{\mathcal{T}} \sum_{l=1}^p X_0^{(l)}(t) \gamma_{1,g}(t) dt, \quad g = 1, \dots, k-1,$$

where $\gamma_{0,g}$ is an intercept term and $\gamma_{1,g}(t)$ is the coefficient function of the predictor. Then, the label G_0 is the one that maximizes $\mathbb{P}(G_0 = g | \mathbf{X}_0)$.

Functional discriminant analysis is another popular method. Similar to the approach used in the multivariate case, the discriminant analysis model (Galeano et al. 2015, Chamroukhi & Nguyen 2019) assumes a set of prior probabilities of belonging to each class, π_i , where $\sum_{i=1}^k \pi_i = 1$. Given the conditional density $f_j(\mathbf{X}) = \mathbb{P}(G = j | \mathbf{X})$, the probability of a new object belonging to class g is

$$\mathbb{P}(G_0 = g | \mathbf{X}_0) = \frac{\pi_g f_g(\mathbf{X}_0)}{\sum_{i=1}^k \pi_i f_i(\mathbf{X}_0)}.$$

One limitation of these methods is that a model is assumed. The regression-type classification assumes a linear model, and the discriminant-type classification assumes a distributional assumption. The inference for these models could be affected if there are outliers in the data. Depth-based classification methods (Hubert et al. 2017) are more useful if the presence of an outlier is suspected.

There are alternative approaches to performing nonparametric classification (Ferraty & Vieu, 2006). A common criterion is the maximum depth procedure (Sguera et al., 2014;

[Cuevas et al., 2007](#)). Without loss of generality, we will assume two groups ($k = 2$), then $\{\mathbf{X}_1^1, \dots, \mathbf{X}_{n_1}^1\}$ and $\{\mathbf{X}_1^2, \dots, \mathbf{X}_{n_2}^2\}$ are the observed curves within each class. For a new observation \mathbf{X}_0 , we assign class 1 if $P_1 > P_2$, where P_i is the depth of observation \mathbf{X}_0 in the sample $\{\mathbf{X}_1^i, \dots, \mathbf{X}_{n_i}^i\}$. This rule can be generalized to the case with k ($k > 2$) groups by identifying the depth of new observation \mathbf{X}_0 in each group and assigning the group with the highest depth to \mathbf{X}_0 . This method can be used with any functional depth.

For multivariate functional data, [Sguera et al. \(2014\)](#) applied a kernel functional spatial depth for supervised classification. [Dai & Genton \(2018c\)](#) investigated supervised classification using directional outlyingness and the outlyingness matrix. [Blanquero et al. \(2019\)](#) have proposed that the support vector machine (SVM) can optimally select the most informative time instants to obtain optimal classification rates. [De Micheaux et al. \(2021\)](#) proposed a curve depth and built a DD-plot to classify unparameterized curves. The cluster and outlier recognition algorithm, introduced by [Qu et al. \(2023\)](#), can also be available in functional classification. The core idea is to compare any functional distance between \mathbf{X}_0 and the center in each group and to determine whether to assign \mathbf{X}_0 to any group based on its rank in the set of distances between curves within the group and the group center. If the distance is larger than a threshold quantile of the set of distances between curves within the group and the group center, then \mathbf{X}_0 is assumed to be an outlier.

6 Discussion

EFDA has broadened its scope from analyzing only fully observed univariate functional data to encompassing irregular multivariate functional data. By using functional depths and distances, EFDA offers a wide array of tools for visualizing, clustering, and classifying both

dense and sparse multivariate functional data, and for detecting outliers.

Functional depths play a pivotal role in establishing functional rankings, forming the foundation for generating functional boxplots and identifying outliers. These functional depths can be categorized into four types, as outlined in [Zuo & Serfling \(2000\)](#). The first type gauges the average closeness of the curve to random samples, exemplified by the BD and the SBD. The second type measures the distance of the curve from random samples, represented by the L^p depth. The third type assesses the outlyingness of a point with respect to the center of random samples, such as the elastic depth. As for the fourth type, it is an index related to the relative depth with respect to the center of the distribution, known as the extremal depth.

Applying functional depths to sparse functional data is more complex because of the irregular coordinate grids. One approach to address this involves estimating curves and their confidence bands or, alternatively, applying global functional depth to sparse functional data directly as proposed by [Qu et al. \(2022\)](#). The original functional boxplot serves to identify central tendencies and outliers. However, it has limitations with regard to detecting certain shape outliers, and it may not be directly applicable to functional data with missing values. To address these limitations, variations of functional boxplots and other visualization tools have been proposed, enabling the detection of shape outliers and facilitating the application of sparse functional data.

Moreover, functional distances play a crucial role in facilitating functional clustering. Examples of such distances include the L^p distance, the total variation distance, and the elastic time distance. By leveraging functional distances, a wide range of classical clustering algorithms, as well as novel ones, can be applied to functional data. To handle common noise present in real-world data, robust two-layer partition clustering techniques can effectively

separate potential outliers from the clusters. In the context of functional classification, both functional depths and distances find utility. When comparing the distance or depth of a curve to all other groups, the group with the smallest distribution proportion of distances (or the maximal depth) is assigned as the label for that curve.

Although this review has extended its domain from classical functional data to sparse multivariate functional data, a wider area of functional data can be considered, and this may pose new challenges in visualization, robust statistics, and clustering/classification. New-generation functional data can include interval-valued functional data ([Nasirzadeh et al. 2022](#); see the simultaneous systolic and diastolic blood pressure of subjects at different visit times), longitudinal functional data from a clinical trial (see the medical imaging data of patients at different time points during a clinical study in the papers by [Adeli et al. 2019](#) and [Zhu et al. 2021](#)), spatial functional data ([Delicado et al. 2010](#); see the longitudinal climate data from arrays of monitors in the nearby area), and wearable health data (see [Smets et al. 2018](#)).

References

- Abraham, C., Cornillon, P. A., Matzner-Løber, E., & Molinari, N. (2003). Unsupervised curve clustering using b-splines. *Scandinavian Journal of Statistics*, 30(3), 581–595.
- Abramowicz, K., Arnqvist, P., Secchi, P., Luna, S. S. d., Vantini, S., & Vitelli, V. (2017). Clustering misaligned dependent curves applied to varved lake sediment for climate reconstruction. *Stochastic Environmental Research and Risk Assessment*, 31(1), 71–85.
- Adeli, E., Meng, Y., Li, G., Lin, W., & Shen, D. (2019). Multi-task prediction of infant cognitive scores from longitudinal incomplete neuroimaging data. *NeuroImage*, 185, 783–792.
- Allen, T. T. (2019). *Introduction to engineering statistics and lean six sigma: statistical quality control and design of experiments and systems*, volume 3. Springer.
- Alvarez-Esteban, P. C., Euán, C., & Ortega, J. (2016). Time series clustering using the total variation distance with applications in oceanography. *Environmetrics*, 27(6), 355–369. env.2398.
- Arribas-Gil, A. & Romo, J. (2014). Shape outlier detection and visualization for functional data: the outliergram. *Biostatistics*, 15(4), 603–619.
- Athanasiadis, S. & Mrkvicka, T. (2019). European insurance market analysis: A multivariate clustering approach. In *Proceedings of the 12th International Scientific Conference INPROFORUM* (pp. 328–333).
- Berrendero, J. R., Justel, A., & Svarc, M. (2011). Principal components for multivariate functional data. *Computational Statistics & Data Analysis*, 55(9), 2619–2634.

- Blanquero, R., Carrizosa, E., Jiménez-Cordero, A., & Martín-Barragán, B. (2019). Variable selection in classification for multivariate functional data. *Information Sciences*, 481, 445–462.
- Bouveyron, C. & Brunet-Saumard, C. (2014). Model-based clustering of high-dimensional data: A review. *Computational Statistics & Data Analysis*, 71, 52 – 78.
- Cai, T. & Yuan, M. (2010). Nonparametric covariance function estimation for functional and longitudinal data. *Technical report, University of Pennsylvania and Georgia institute of technology*.
- Carroll, C., Müller, H.-G., & Kneip, A. (2021). Cross-component registration for multivariate functional data, with application to growth curves. *Biometrics*, 77(3), 839–851.
- Castruccio, S., Genton, M. G., & Sun, Y. (2019). Visualizing spatiotemporal models with virtual reality: from fully immersive environments to applications in stereoscopic view. *Journal of the Royal Statistical Society: Series A*, 182(2), 379–387.
- Chamroukhi, F. & Nguyen, H. D. (2019). Model-based clustering and classification of functional data. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(4), e1298.
- Chang, W., Cheng, J., Allaire, J., Xie, Y., & McPherson, J. (2015). Package ‘shiny’. <https://cran.r-project.org/web/packages/shiny/shiny.pdf>.
- Claeskens, G., Hubert, M., Slaets, L., & Vakili, K. (2014). Multivariate functional halfspace depth. *Journal of the American Statistical Association*, 109(505), 411–423.
- Cuesta-Albertos, J. A. & Fraiman, R. (2007). Impartial trimmed K-means for functional data. *Computational Statistics & Data Analysis*, 51(10), 4864–4877.

- Cuevas, A., Febrero, M., & Fraiman, R. (2007). Robust estimation and classification for functional data via projection-based depth notions. *Computational Statistics*, 22(3), 481–496.
- Dai, W., Athanasiadis, S., & Mrkvička, T. (2021). A new functional clustering method with combined dissimilarity sources and graphical interpretation. In *Computational Statistics and Applications* (pp. 3–23). IntechOpen.
- Dai, W. & Genton, M. G. (2018a). Functional boxplots for multivariate curves. *Stat*, 7(1), e190.
- Dai, W. & Genton, M. G. (2018b). Multivariate functional data visualization and outlier detection. *Journal of Computational and Graphical Statistics*, 27(4), 923–934.
- Dai, W. & Genton, M. G. (2018c). An outlyingness matrix for multivariate functional data classification. *Statistica Sinica*, 28(4), 2435–2454.
- Dai, W. & Genton, M. G. (2019). Directional outlyingness for multivariate functional data. *Computational Statistics & Data Analysis*, 131, 50–65.
- De Micheaux, P. L., Mozharovskyi, P., & Vimond, M. (2021). Depth for curve data and applications. *Journal of the American Statistical Association*, 116(536), 1881–1897.
- Delicado, P., Giraldo, R., Comas, C., & Mateu, J. (2010). Statistics for spatial functional data: some recent contributions. *Environmetrics: The official journal of the International Environmetrics Society*, 21(3-4), 224–239.
- Euán, C., Ombao, H., & Ortega, J. (2018). The hierarchical spectral merger algorithm: A new time series clustering procedure. *Journal of Classification*, 35(1), 71 – 99.

- Euán, C. & Sun, Y. (2019). Directional spectra-based clustering for visualizing patterns of ocean waves and winds. *Journal of Computational and Graphical Statistics*, 28(3), 659–670.
- Euán, C., Sun, Y., & Ombao, H. (2019). Coherence-based time series clustering for statistical inference and visualization of brain connectivity. *Ann. Appl. Statist.*, 13(2), 990–1015.
- Fan, J. & Gijbels, I. (1996). *Local polynomial modelling and its applications: monographs on statistics and applied probability 66*, volume 66. CRC Press.
- Ferraty, F. & Vieu, P. (2006). *Nonparametric Functional Data Analysis: Theory and Practice*. Springer.
- Fraiman, R. & Muniz, G. (2001). Trimmed means for functional data. *TEST*, 10, 419–440.
- Friedman, J. H. & Stuetzle, W. (2002). John W. Tukey’s work on interactive graphics. *The Annals of Statistics*, 30(6), 1629–1639.
- Galeano, P., Joseph, E., & Lillo, R. E. (2015). The Mahalanobis distance for functional data with applications to classification. *Technometrics*, 57(2), 281–291.
- Genton, M. G., Castruccio, S., Crippa, P., Dutta, S., Huser, R., Sun, Y., & Vettori, S. (2015). Visuanimation in statistics. *Stat*, 4(1), 81–96.
- Genton, M. G., Johnson, C., Potter, K., Stenchikov, G., & Sun, Y. (2014). Surface boxplots. *Stat*, 3(1), 1–11.
- Happ, C. & Greven, S. (2018). Multivariate functional principal component analysis for data observed on different (dimensional) domains. *Journal of the American Statistical Association*, 113(522), 649–659.

- Harris, T., Tucker, J. D., Li, B., & Shand, L. (2021). Elastic depths for detecting shape anomalies in functional data. *Technometrics*, 63(4), 466–476.
- Heckman, N. & Zamar, R. (2000). Comparing the shapes of regression functions. *Biometrika*, 87(1), 135–144.
- Hitchcock, D. B. & Greenwood, M. C. (2015). Clustering functional data. In *Handbook of Cluster Analysis* chapter 13, (pp. 265–287). Chapman and Hall/CRC.
- Hong, Y., Davis, B., Marron, J., Kwitt, R., Singh, N., Kimbell, J. S., Pitkin, E., Superfine, R., Davis, S. D., Zdanski, C. J., et al. (2014). Statistical atlas construction via weighted functional boxplots. *Medical image analysis*, 18(4), 684–698.
- Horváth, L. & Kokoszka, P. (2012). *Inference for Functional Data with Applications*. Springer.
- Hsing, T. & Eubank, R. (2015). *Theoretical Foundations of Functional Data Analysis, with an Introduction to Linear Operators*, volume 997. John Wiley & Sons.
- Huang, H. & Sun, Y. (2019). Visualization and assessment of spatio-temporal covariance properties. *Spatial Statistics*, 34, 100272.
- Hubert, M., Rousseeuw, P., & Segaert, P. (2017). Multivariate and functional classification using depth and distance. *Advances in Data Analysis and Classification*, 11, 445–466.
- Hubert, M., Rousseeuw, P. J., & Segaert, P. (2015). Multivariate functional outlier detection. *Statistical Methods & Applications*, 24(2), 177–202.
- Hyndman, R. J. (1996). Computing and graphing highest density regions. *The American Statistician*, 50(2), 120–126.

- Hyndman, R. J. & Shang, H. L. (2010). Rainbow plots, bagplots, and boxplots for functional data. *Journal of Computational and Graphical Statistics*, 19(1), 29–45.
- Ieva, F. & Paganoni, A. M. (2013). Depth measures for multivariate functional data. *Communications in Statistics-Theory and Methods*, 42(7), 1265–1276.
- Ieva, F., Paganoni, A. M., Pigoli, D., & Vitelli, V. (2013). Multivariate functional clustering for the morphological analysis of electrocardiograph curves. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 62(3), 401–418.
- Jacques, J. & Preda, C. (2014). Functional data clustering: a survey. *Advances in Data Analysis and Classification*, 8(3), 231–255.
- Jörnsten, R. (2004). Clustering and classification based on the L1 data depth. *Journal of Multivariate Analysis*, 90(1), 67–89.
- Kuhnt, S. & Rehage, A. (2016). An angle-based multivariate functional pseudo-depth for shape outlier detection. *Journal of Multivariate Analysis*, 146, 325–340.
- Li, C., Xiao, L., & Luo, S. (2020). Fast covariance estimation for multivariate sparse functional data. *Stat*, 9(1), e245.
- López-Pintado, S. & Romo, J. (2009). On the concept of depth for functional data. *Journal of the American Statistical Association*, 104(486), 718–734.
- López-Pintado, S. & Romo, J. (2011). A half-region depth for functional data. *Computational Statistics & Data Analysis*, 55(4), 1679–1695.
- López-Pintado, S., Sun, Y., Lin, J. K., & Genton, M. G. (2014). Simplicial band depth for multivariate functional data. *Advances in Data Analysis and Classification*, 8, 321–338.

Lòpez-Pintado, S. & Wei, Y. (2011). Depth for sparse functional data. In *Recent Advances in Functional Data Analysis and Related Topics* (pp. 209–212).: Springer.

Martinez, W. L., Martinez, A. R., & Solka, J. L. (2017). *Exploratory data analysis with MATLAB®*. Chapman and Hall/CRC.

Mirzargar, M., Whitaker, R. T., & Kirby, R. M. (2014). Curve boxplot: Generalization of boxplot for ensembles of curves. *IEEE Transactions on Visualization and Computer Graphics*, 20(12), 2654–2663.

Myllymäki, M. & Mrkvíčka, T. (2019). GET: Global envelopes in R. *arXiv preprint arXiv:1911.06583*.

Myllymäki, M., Mrkvíčka, T., Grabarnik, P., Seijo, H., & Hahn, U. (2017). Global envelope tests for spatial processes. *Journal of the Royal Statistical Society: Series B*, 79(2), 381–404.

Nagy, S., Gijbels, I., & Hlubinka, D. (2017). Depth-based recognition of shape outlying functions. *Journal of Computational and Graphical Statistics*, 26(4), 883–893.

Narisetty, N. N. & Nair, V. N. (2016). Extremal depth for functional data and applications. *Journal of the American Statistical Association*, 111(516), 1705–1714.

Nasirzadeh, R., Nasirzadeh, F., & Mohammadi, Z. (2022). Some non-parametric regression models for interval-valued functional data. *Stat*, 11(1), e443.

Ojo, O. T., Fernández Anta, A., Genton, M. G., & Lillo, R. E. (2023). Multivariate functional outlier detection using the fast massive unsupervised outlier detection indices. *Stat*, 12(1), e567.

- Ojo, O. T., Fernández Anta, A., Lillo, R. E., & Sguera, C. (2022). Detecting and classifying outliers in big functional data. *Advances in Data Analysis and Classification*, 16(3), 725–760.
- Park, H.-S. & Jun, C.-H. (2009). A simple and fast algorithm for k-medoids clustering. *Expert systems with applications*, 36(2), 3336–3341.
- Qu, Z., Dai, W., & Genton, M. G. (2022). Global depths for irregularly observed multivariate functional data. *arXiv preprint arXiv:2211.15125*.
- Qu, Z., Dai, W., & Genton, M. G. (2023). Robust two-layer partition clustering of sparse multivariate functional data. *Econometrics and Statistics*, In press.
- Qu, Z. & Genton, M. G. (2022). Sparse functional boxplots for multivariate curves. *Journal of Computational and Graphical Statistics*, 31(4), 976–989.
- Raj, M., Mirzargar, M., Ricci, R., Kirby, R. M., & Whitaker, R. T. (2017). Path boxplots: a method for characterizing uncertainty in path ensembles on a graph. *Journal of Computational and Graphical Statistics*, 26(2), 243–252.
- Ramsay, J., Hooker, G., & Graves, S. (2009). Introduction to functional data analysis. In *Functional data analysis with R and MATLAB* (pp. 1–19). Springer.
- Ramsay, J. O. & Dalzell, C. (1991). Some tools for functional data analysis. *Journal of the Royal Statistical Society: Series B (Methodological)*, 53(3), 539–561.
- Ramsay, J. O. & Silverman, B. W. (2005). *Functional Data Analysis* (second ed.). Springer.
- Rivera-García, D., García-Escudero, L. A., Mayo-Iscar, A., & Ortega, J. (2019). Robust clustering for functional data based on trimming and constraints. *Advances in Data Analysis and Classification*, 13(1), 201–225.

- Rousseeuw, P. J., Raymaekers, J., & Hubert, M. (2018). A measure of directional outlyingness with applications to image data and video. *Journal of Computational and Graphical Statistics*, 27(2), 345–359.
- Rousseeuw, P. J., Ruts, I., & Tukey, J. W. (1999). The bagplot: a bivariate boxplot. *The American Statistician*, 53(4), 382–387.
- Sangalli, L. M., Secchi, P., Vantini, S., & Veneziani, A. (2009). A case study in exploratory functional data analysis: geometrical features of the internal carotid artery. *Journal of the American Statistical Association*, 104(485), 37–48.
- Sangalli, L. M., Secchi, P., Vantini, S., & Vitelli, V. (2010). K-mean alignment for curve clustering. *Computational Statistics & Data Analysis*, 54(5), 1219 – 1233.
- Schmutz, A., Jacques, J., Bouveyron, C., Cheze, L., & Martin, P. (2020). Clustering multivariate functional data in group-specific functional subspaces. *Computational Statistics*, 35(3), 1101–1131.
- Serban, N. & Wasserman, L. (2005). Cats: Clustering after transformation and smoothing. *Journal of the American Statistical Association*, 100(471), 990–999.
- Serfling, R. & Wijesuriya, U. (2017). Depth-based nonparametric description of functional data, with emphasis on use of spatial depth. *Computational Statistics & Data Analysis*, 105, 24–45.
- Sguera, C., Galeano, P., & Lillo, R. (2014). Spatial depth-based classification for functional data. *TEST*, 23(4), 725–750.
- Sguera, C. & López-Pintado, S. (2021). A notion of depth for sparse functional data. *TEST*, 30(3), 630–649.

- Shang, H. L. (2014). A survey of functional principal component analysis. *AStA Advances in Statistical Analysis*, 98, 121–142.
- Shang, H. L. & Hyndman, R. J. (2019). Package ‘rainbow’. *R packages*, 43.
- Sievert, C., Hocking, T., Chamberlain, S., Ram, K., Corvellec, M., & Despouy, P. (2018). plotly for R.
- Smets, E., Rios Velazquez, E., Schiavone, G., Chakroun, I., D’Hondt, E., De Raedt, W., Cornelis, J., Janssens, O., Van Hoecke, S., Claes, S., et al. (2018). Large-scale wearable data reveal digital phenotypes for daily-life stress detection. *NPJ digital medicine*, 1(1), 67.
- Srivastava, A., Wu, W., Kurtek, S., Klassen, E., & Marron, J. S. (2011). Registration of functional data using Fisher-Rao metric. *arXiv preprint arXiv:1103.3817*.
- Sun, Y. & Genton, M. G. (2011). Functional boxplots. *Journal of Computational and Graphical Statistics*, 20(2), 316–334.
- Sun, Y. & Genton, M. G. (2012). Adjusted functional boxplots for spatio-temporal data visualization and outlier detection. *Environmetrics*, 23(1), 54–64.
- Tukey, J. W. (1975). Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians*, volume 2 (pp. 523–531).
- Tukey, J. W. (1977). *Exploratory data analysis*, volume 2. Addison-Wesley. Reading, MA, USA.
- Tukey, J. W. (1980). We need both exploratory and confirmatory. *The American Statistician*, 34(1), 23–25.
- Wand, M. P. & Jones, M. C. (1995). *Kernel smoothing*. London: Chapman and Hall.

Wang, J. L., Chiou, J. M., & Müller, H. G. (2016). Functional data analysis. *Annual Review of Statistics and Its Application*, 3, 257–295.

Wang, Y. (2011). *Smoothing Splines: Methods and Applications*. New York: Chapman and Hall.

Whitaker, R. T., Mirzargar, M., & Kirby, R. M. (2013). Contour boxplots: A method for characterizing uncertainty in feature sets from simulation ensembles. *IEEE Transactions on Visualization and Computer Graphics*, 19(12), 2713–2722.

Wrobel, J. & Goldsmith, J. (2016). refund. shiny: Interactive plotting for functional data analyses. *R package version 0.3. 0, URL https://CRAN.R-project.org/package=refund.shiny*.

Wrobel, J., Park, S. Y., Staicu, A. M., & Goldsmith, J. (2016). Interactive graphics for functional data analyses. *Stat*, 5(1), 108–118.

Xiao, L., Li, C., Checkley, W., & Crainiceanu, C. (2018). Fast covariance estimation for sparse functional data. *Statistics and computing*, 28, 511–522.

Xie, W., Kurtek, S., Bharath, K., & Sun, Y. (2017). A geometric approach to visualization of variability in functional data. *Journal of the American Statistical Association*, 112(519), 979–993.

Yao, F., Müller, H. G., & Wang, J. L. (2005). Functional linear regression analysis for longitudinal data. *The Annals of Statistics*, 33(6), 2873–2903.

Yao, Z., Dai, W., & G. Genton, M. (2020). Trajectory functional boxplots. *Stat*, 9(1), e289.

Zhao, X., Marron, J., & Wells, M. T. (2004). The functional data analysis view of longitudinal data. *Statistica Sinica*, 14(2004), 789–808.

Zhu, Y., Kim, M., Zhu, X., Kaufer, D., Wu, G., Initiative, A. D. N., et al. (2021). Long range early diagnosis of Alzheimer's disease using longitudinal MR imaging data. *Medical image analysis*, 67, 101825.

Zuo, Y. & Serfling, R. (2000). General notions of statistical depth function. *The Annals of Statistics*, 28(2), 461–482.