

Sparse Functional Boxplots for Multivariate Curves

Zhuo Qu, Marc G. Genton

Statistics Program, King Abdullah University of Science and Technology, Saudi Arabia

October 26, 2022

Overview: Sparse Functional Boxplots for Multivariate Curves

1 Motivation

- Restrictions of current visualization tools: measured on common grids
- Limited discussion of depths for sparse multivariate functional data

2 Preparatory Work

- Fitting sparse multivariate functional data
- Depths for sparse multivariate functional data

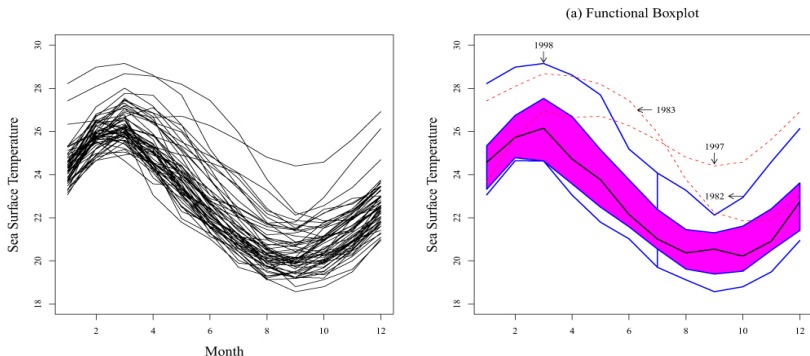
3 Construction of Sparse Visualization Tools

- Sparse functional boxplot
- Intensity sparse functional boxplot

4 Application: Malnutrition Data

1. Motivation I: Restriction of the current visualization tool

- ① **Functional Boxplot** (Sun & Genton 2011) is an exploratory visualization tool for functional data



- ② **Enhanced** or **adjusted** functional boxplots (Sun & Genton 2012), **surface** boxplot (Genton et al. 2014), **two-stage** functional boxplot (Dai & Genton 2018), **trajectory** functional boxplot (Yao et al. 2020)

Motivation: Malnutrition Example

How about data with missing values? (Apply them to our tool: Figure 5)

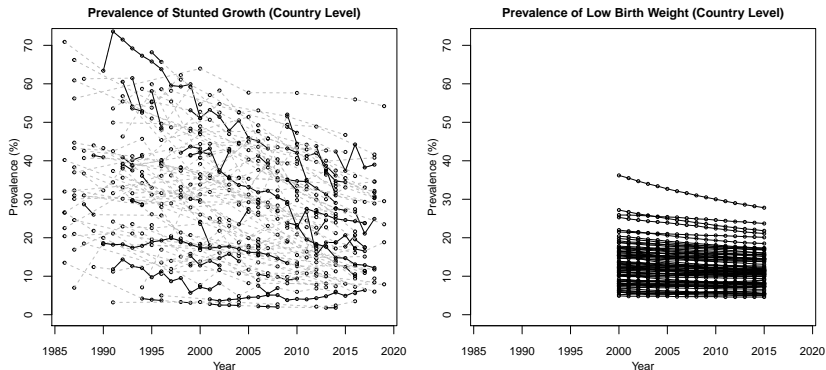


Figure: The observed prevalence of stunted growth and prevalence of low birth weight for 77 countries from 1985 to 2019 (source). Observations are joined with solid black lines if observed continuously; otherwise, joined with gray dashed lines.

Motivation II: Limited discussion of ordering sparse multivariate functional data

- Visualization tool for sparse multivariate functional data



Ranking sparse multivariate functional data

- Current methods to order multivariate functional data:
 - 1 **Depth**-based: weighted average of marginal functional depth (WMFD, Ieva & Paganoni 2013), modified simplicial band depth (MSBD, López-Pintado et al. 2014), multivariate functional halfspace depth (MFHD, Claeskens et al. 2014)
 - 2 **Outlying**-based: directional outlyingness (DO, Dai & Genton 2019)

Restrictions:

- ① Visualization tools for univariate (multivariate) functional data
- ② Methods for ordering multivariate functional data



assume that data are observed on common grids

Thread:

- Depths for sparse multivariate functional data (✗, not considered)
- Depths for univariate sparse functional data (✓, López-Pintado & Wei (2011) and Sguera & López-Pintado (2020))

Logic:

- ① Impute missing values
- ② Propose multivariate functional depth
- ③ Add features of sparseness in functional boxplot

Contributions:

- ① improve sparse multivariate functional data fitting
- ② consider possible depths for sparse multivariate functional data
- ③ propose exploratory visualization tools for both univariate and multivariate functional data with missing values

2. Preparatory Work

To propose the depth for sparse multivariate functional data, recall:

I. Happ & Greven (2018) proposed the **principal component analysis for multivariate functional data (MFPCA)** that are defined on different time domains



I. explore the improvement of sparse multivariate data fitting

II. López-Pintado & Wei (2011) and Sguera & López-Pintado (2020) proposed possible **depths for univariate sparse functional data**



II. consider several depths to order sparse multivariate functional data

Preparatory Work I: Data Fitting

Data Notation:

- $\mathcal{T} = \mathcal{T}_1 \times \cdots \times \mathcal{T}_p$, where $(\mathcal{T}_j)_{j=1}^p$ is a compact set in \mathbb{R}^{d_j}
- $\mathbf{Y}(\mathbf{t}) = (Y^{(1)}(t^{(1)}), \dots, Y^{(p)}(t^{(p)}))^\top$, $\mathbf{t} := (t^{(1)}, \dots, t^{(p)}) \in \mathcal{T}$
- $Y^{(j)}(t^{(j)}) : \mathcal{T}_j \rightarrow \mathbb{R}$ is assumed to be square-integrable in \mathcal{T}_j
- $\tilde{\mathbf{Y}}(\mathbf{t}) = (\tilde{Y}^{(1)}(t^{(1)}), \dots, \tilde{Y}^{(p)}(t^{(p)}))^\top \in \mathbb{R}^p$, due to measurement errors $\boldsymbol{\epsilon}_i = (\epsilon_i^{(1)}, \dots, \epsilon_i^{(p)})^\top$ with $\epsilon_i^{(j)} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_j^2)$

The i th ($i = 1, \dots, N$) observation $\tilde{\mathbf{Y}}_i(\mathbf{t}_i)$ is denoted as

$$\tilde{\mathbf{Y}}_i(\mathbf{t}_i) = \mathbf{Y}_i(\mathbf{t}_i) + \boldsymbol{\epsilon}_i = \boldsymbol{\mu}(\mathbf{t}_i) + \sum_{m=1}^{\infty} \rho_{i,m} \boldsymbol{\psi}_{i,m}(\mathbf{t}_i) + \boldsymbol{\epsilon}_i,$$

where $\boldsymbol{\mu}(\mathbf{t}_i)$ is the mean function, $\boldsymbol{\psi}_{i,m}(\mathbf{t}_i)$ and $\rho_{i,m}$ are the multivariate eigenfunctions and eigenscores respectively

Preparatory Work I: Data Fitting

Happ & Greven (2018) proposed the MFPCA fit based on unobserved MFPC decomposition objects $\theta = \{\mu, \rho, \psi\}$:

$$\hat{Y}_{\hat{\theta},i} := E[\tilde{Y}_i | \hat{\theta}] = \hat{\mu}_i + \sum_{m=1}^M \hat{\rho}_{i,m} \hat{\psi}_m, \quad i = 1, \dots, N$$

- 1 Due to the uncertainty of the obtained eigenvalues and eigenfunctions (Goldsmith et al, 2013), we propose the **bootstrap improved MFPCA fit (BMFPCA fit)**:

$$\hat{Y}_i = E_{\hat{\theta}} \left\{ E_{\tilde{Y}_i | \hat{\theta}} [\tilde{Y}_i | \hat{\theta}] \right\}$$

- 2 **$(1 - \alpha)$ bootstrap confidence band**: we take the $(1 - \frac{\alpha}{2})$ th and the $\frac{\alpha}{2}$ th percentiles across the bootstrap as the confidence upper bound $\hat{Y}_{ub,i}$ and lower bound $\hat{Y}_{lb,i}$, respectively

Preparatory Work II: Possible Depth Notions

Take multivariate functional halfspace depth (MFHD) as an example

$$MFHD(\mathbf{X}; F_{\mathbf{Y}}, \beta) = \int_{\mathcal{T}} HD(\mathbf{X}(t); F_{\mathbf{Y}(t)}) \cdot w_{\beta}(t; F_{\mathbf{Y}(t)}) dt$$

Note:

- ① $HD(\mathbf{X}(t); F_{\mathbf{Y}(t)}): \mathbb{R}^p \mapsto \mathbb{R}$ is the halfspace depth of $\mathbf{X}(t)$ with respect to the random variable with cumulative distribution function (cdf) $F_{\mathbf{Y}(t)}$
- ② $w_{\beta}(t; F_{\mathbf{Y}(t)})$ is a weight function satisfying $\int_{\mathcal{T}} w_{\beta}(t; F_{\mathbf{Y}(t)}) dt = 1$ for a fixed $\beta \in (0, 1]$

Preparatory Work II: Possible Depth Notions

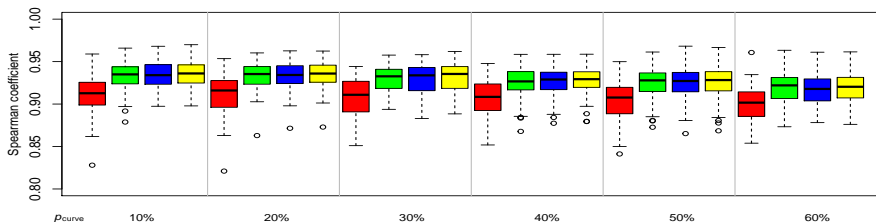
Let $\hat{\mathbf{Y}}^f = \{\hat{\mathbf{Y}}_1, \dots, \hat{\mathbf{Y}}_N\}$ be a set of fitted data, $\hat{\mathbf{Y}}_{ub}^f = \{\hat{\mathbf{Y}}_{ub,1}, \dots, \hat{\mathbf{Y}}_{ub,N}\}$ be a set of confidence upper bounds, $\hat{\mathbf{Y}}_{lb}^f = \{\hat{\mathbf{Y}}_{lb,1}, \dots, \hat{\mathbf{Y}}_{lb,N}\}$ be a set of confidence lower bounds, and $\hat{\mathbf{Y}}_{upd}^f = \{\hat{\mathbf{Y}}^f, \hat{\mathbf{Y}}_{ub}^f, \hat{\mathbf{Y}}_{lb}^f\}$

- 1 Apply MFHD to each element in $\hat{\mathbf{Y}}^f$: $MFHD(\hat{\mathbf{Y}}; \hat{\mathbf{Y}}^f)$

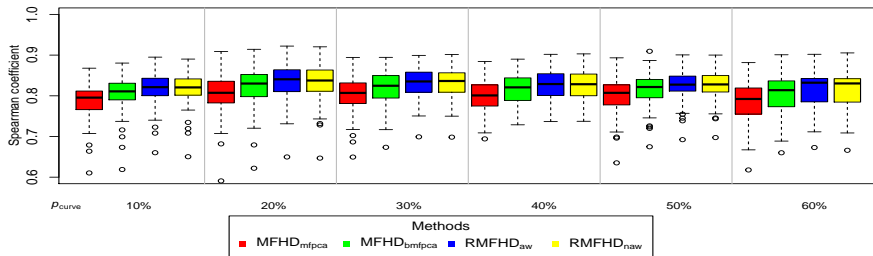
$$\begin{aligned} & \text{RMFHD}_{\text{type}}(\hat{\mathbf{Y}}; \hat{\mathbf{Y}}_{upd}) \\ &= \begin{cases} \frac{1}{3}MFHD(\hat{\mathbf{Y}}; \hat{\mathbf{Y}}_{upd}) + \frac{1}{3}MFHD(\hat{\mathbf{Y}}_{ub}; \hat{\mathbf{Y}}_{upd}) \\ \quad + \frac{1}{3}MFHD(\hat{\mathbf{Y}}_{lb}; \hat{\mathbf{Y}}_{upd}), & \text{if type} = \text{aw}, \\ \frac{1}{2}MFHD(\hat{\mathbf{Y}}; \hat{\mathbf{Y}}_{upd}) + \frac{1}{4}MFHD(\hat{\mathbf{Y}}_{ub}; \hat{\mathbf{Y}}_{upd}) \\ \quad + \frac{1}{4}MFHD(\hat{\mathbf{Y}}_{lb}; \hat{\mathbf{Y}}_{upd}), & \text{if type} = \text{naw}. \end{cases} \end{aligned}$$

Optimal Depth: MFHD_{bmfpc}

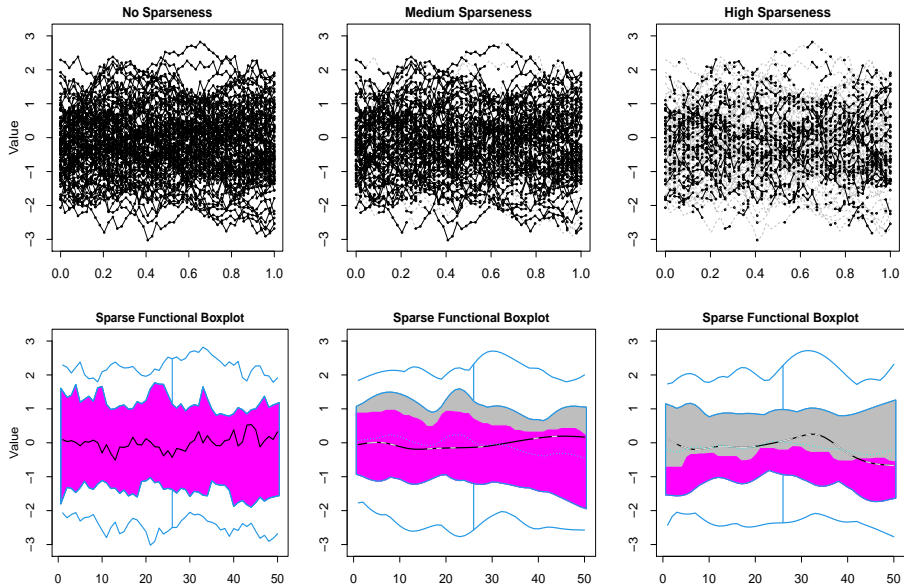
Spearman coefficients from various methods: Model 1 (no outlier case)



Spearman coefficients from various methods: Model 8 (covariance outlier case)



Examples of the Sparse Functional Boxplot



3. Construction of Sparse Visualization Tools

- Let $Z_{[r]}^{(j)}(t)$ be the j th component of the r th deepest curve evaluated at $t \in \mathcal{T}$, and $\lceil n/2 \rceil$ be the smallest integer $\geq n/2$

① **Central region:**

$$C_{0.5}^{(j)} := \{(t, Z^{(j)}(t)) : \min_{r=1, \dots, \lceil n/2 \rceil} Z_{[r]}^{(j)}(t) \leq Z^{(j)}(t) \leq \max_{r=1, \dots, \lceil n/2 \rceil} Z_{[r]}^{(j)}(t)\}$$

② **The median:** $Z_{[1]}^{(j)}(t)$

- ③ **The outliers:** A curve \mathbf{Z}_o is classified in S_o ($\mathbf{Z}_o \in S_o$), if $Z_o^{(j)}(t) > Z_{ub,0.5}^{(j)}(t) + 1.5R_{0.5}^{(j)}(t)$ or $Z_o^{(j)}(t) < Z_{lb,0.5}^{(j)}(t) - 1.5R_{0.5}^{(j)}(t)$

- ④ **The non-outlying maximal (minimal) bound:** the non-outlying maximal bound $Z_{ub}^{(j)}(t) := \max_{\mathbf{Z} \in \mathcal{Z}^f \setminus S_o} Z^{(j)}(t)$

Construction of Sparse Functional Boxplot

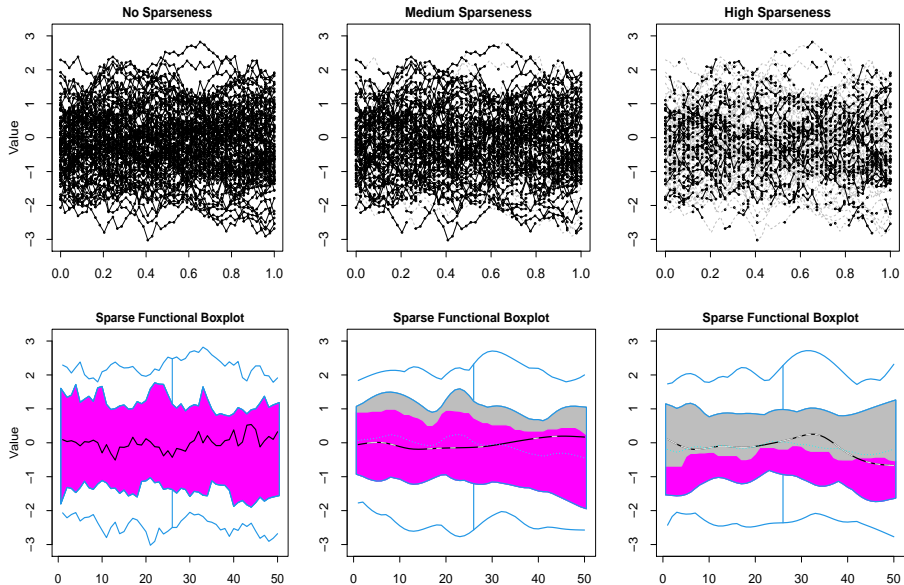
The sparse functional boxplot, apart from displaying the aforementioned features, underlines the sparseness features in **the median** $Z_{[1]}^{(j)}(t)$, **the 50% central region** $C_{0.5}^{(j)}$, and **the detected outliers** S_o

- 1 **The median and outliers:** underline the missing values in gray
- 2 **The central region:**

- At each time point $t \in \mathcal{T}$, the sparseness proportion $p_s^{(j)}(t) := \frac{n_{ms}^{(j)}(t)}{\{n_{ms}^{(j)}(t) + n_{obs}^{(j)}(t)\}}$ within $C_{0.5}^{(j)}$
- Define the proportion line $l^{(j)}(t, p_s^{(j)}(t)) := Z_{lb,0.5}^{(j)}(t) - p_s^{(j)}(t)R_{0.5}^{(j)}(t)$ for $t \in \mathcal{T}$

Show the observed proportion below the proportion line in **magenta** and the sparseness proportion above in **gray**

Examples of the Sparse Functional Boxplot



Construction of Intensity Sparse Functional Boxplot

In addition, we display **the intensity** of fitted missing point patterns within $C_{0.5}^{(j)}$, which is expressed in estimated missing points per unit area

Definition of the Intensity:

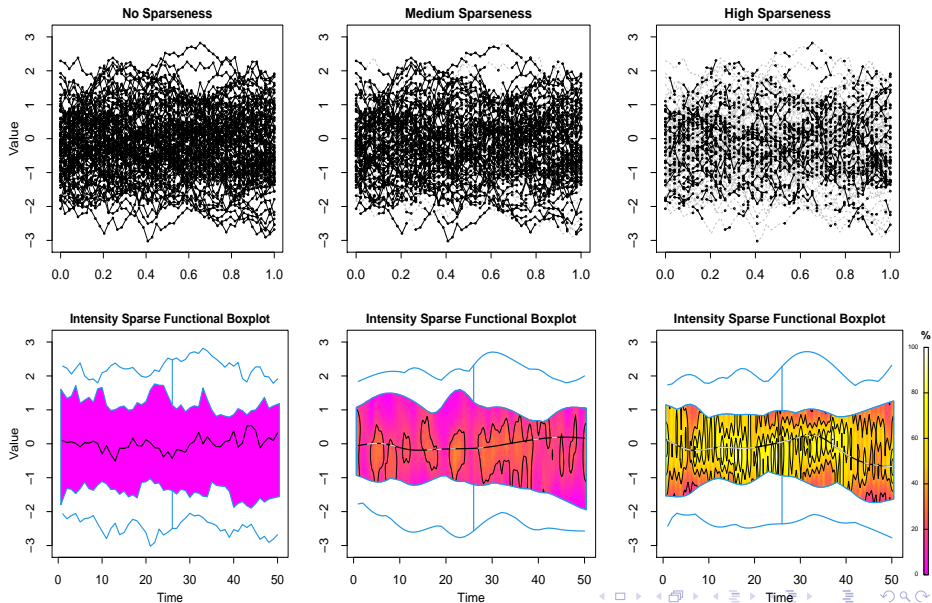
- Regard S fitted sparse points within $C_{0.5}^{(j)}$ as a spatial point pattern $\mathbf{u}_s^{(j)} := \{(t_s, Z_s^{(j)}) \in C_{0.5}^{(j)}, s = 1, \dots, S\}$ with t_s the time, and $Z_s^{(j)}$ the fitted value inside the central region

- The sparseness intensity at a new point $\mathbf{u}^{(j)} \in C_{0.5}^{(j)}$ is

$$\lambda(\mathbf{u}^{(j)}) = e(\mathbf{u}^{(j)}) \sum_{s=1}^S w_s \mathcal{K}(\mathbf{u}_s^{(j)} - \mathbf{u}^{(j)}),$$

where \mathcal{K} is the Gaussian smoothing kernel, $e(\mathbf{u}^{(j)})$ is an edge correction factor, and w_s is the weight

Examples of the Intensity Sparse Functional Boxplot



4. Malnutrition data: Sparse Functional Boxplot

As shown in Figure 1, we have two variables: prevalence of **stunted growth**, and prevalence of **low birth weight** from 77 countries from 1985 to 2019, which belong to the point sparseness and partial sparseness, respectively.

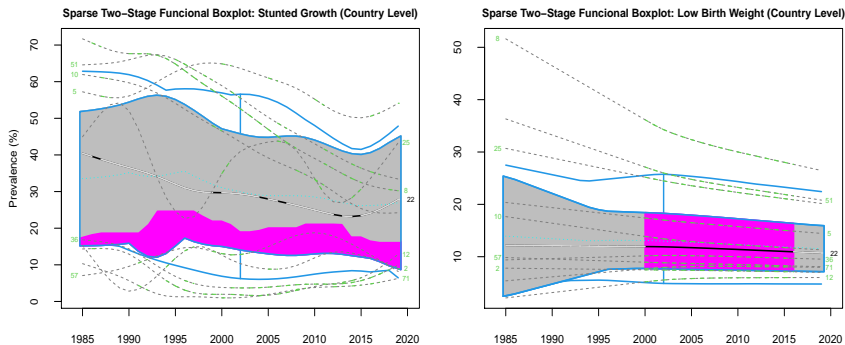


Figure: Visualization of stunted growth and low birth weight data for 77 countries with the sparse two-stage functional boxplot.

Malnutrition data: Intensity Sparse Functional Boxplot

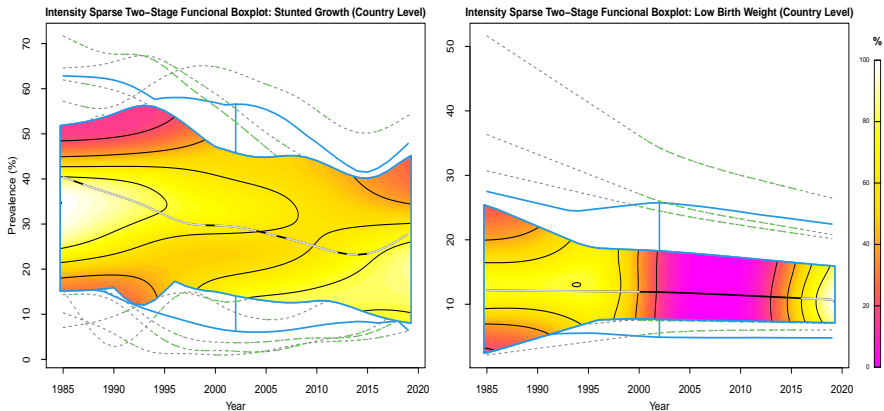


Figure: Visualization of stunted growth and low birth weight data for 77 countries with the intensity sparse two-stage functional boxplot.

- Zhuo Qu & Marc G. Genton (2022) Sparse Functional Boxplots for Multivariate Curves, Journal of Computational and Graphical Statistics, DOI: 10.1080/10618600.2022.2066680
- To replicate the data analysis/apply your data to our visualization tools, please visit: Github

References

- Tukey, John W (1975) Mathematics and the picturing of data *Proceedings of the International Congress of Mathematicians, Vancouver 2*, 525 – 531.
- López-Pintado, Sara and Romo, Juan (2009) On the concept of depth for functional data *Journal of the American statistical Association* 104 (486), 718 – 734.
- Sun, Ying and Genton, Marc G (2011) Functional boxplots *Journal of Computational and Graphical Statistics* 20 (2), 316 – 334.
- Sun, Ying and Genton, Marc G (2012) Adjusted functional boxplots for spatio-temporal data visualization and outlier detection *Environmetrics* 23 (1), 54 – 64.
- Genton, Marc G and Johnson, Christopher and Potter, Kristin and Stenchikov, Georgiy and Sun, Ying (2014) Surface boxplots *Stat* 3 (1), 1 – 11.
- Dai, Wenlin and Genton, Marc G (2018) Functional boxplots for multivariate curves *Stat* 7 (1), e190.
- Yao, Zonghui and Dai, Wenlin and G. Genton, Marc (2020) Trajectory functional boxplots *Stat* 9 (1), e289.
- Qu, Zhuo and Genton, Marc G (2021) Sparse Functional Boxplots for Multivariate Curves *arXiv preprint arXiv:2103.07868*
- Ieva, Francesca and Paganoni, Anna M (2013) Depth measures for multivariate functional data *Communications in Statistics-Theory and Methods* 42 (7), 1265 – 1276.

References

- López-Pintado, Sara and Sun, Ying and Lin, Juan K and Genton, Marc G (2014) Simplicial band depth for multivariate functional data *Advances in Data Analysis and Classification* 8 (3), 321 – 338.
- Claeskens, Gerda and Hubert, Mia and Slaets, Leen and Vakili, Kaveh (2014) Multivariate functional halfspace depth *Journal of the American Statistical Association* 109 (505), 411 – 423.
- Dai, Wenlin and Genton, Marc G (2019) Directional outlyingness for multivariate functional data *Computational Statistics & Data Analysis* 131, 50 – 65.
- López-Pintado, Sara and Wei, Ying (2011) Depth for sparse functional data *Recent Advances in Functional Data Analysis and Related Topics* 209 – 212.
- Sguera, Carlo and López-Pintado, Sara (2020) A notion of depth for sparse functional data *TEST* 1 – 20.
- Happ, Clara and Greven, Sonja (2018) Multivariate functional principal component analysis for data observed on different (dimensional) domains *Journal of the American Statistical Association* 113 (522) 649 – 659.
- Goldsmith, Jeff and Greven, Sonja and Crainiceanu, CIPRIAN (2013) Corrected confidence bands for functional data using principal components *Biometrics* 69 (1) 41 – 51.