



**HAL**  
open science

# Regression imputation in the functional linear model with missing values in the response

Christophe Crambes, Yousri Henchiri

► **To cite this version:**

Christophe Crambes, Yousri Henchiri. Regression imputation in the functional linear model with missing values in the response. 2017. hal-01521954v3

**HAL Id: hal-01521954**

**<https://hal.science/hal-01521954v3>**

Preprint submitted on 3 Oct 2017 (v3), last revised 11 Mar 2019 (v5)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Regression imputation in the functional linear model with missing values in the response

Christophe Crambes · Yousri Henchiri

Received: date / Accepted: date

**Abstract** We are interested in functional linear regression when some observations of the real response are missing, while the functional covariate is completely observed. A complete case regression imputation method of missing data is presented, using functional principal component regression to estimate the functional coefficient of the model. We study the asymptotic behaviour of the error when the missing data are replaced by the regression imputed value, in a 'missing at random' framework. The completed database can be used to estimate the functional coefficient of the model and to predict new values of the response. The practical behaviour of the method is also studied on simulated data sets. A real dataset study is performed in the environmental context of air quality.

**Keywords** Functional linear model; Missing data; Missing at random; Principal components regression; Mean square error prediction

**Mathematics Subject Classification (2000)** 62G20 · 62J05 · 62F12 · 62P12

## 1 Introduction

Literature on functional data is really wide, as attested by the numerous books on this subject these last years. The estimation and forecasting theories of lin-

---

Christophe Crambes  
Institut Montpellierain Alexander Grothendieck (IMAG), Université de Montpellier, France  
E-mail: crambes.christophe@univ-montp2.fr

Yousri Henchiri  
Université de la Manouba, Institut Supérieur des Arts Multimédia de la Manouba (ISAMM),  
Tunisie.  
Université de Tunis El Manar, Laboratoire de Modélisation Mathématique et Numérique  
dans les Sciences de l'Ingénieur (ENIT-LAMSIN), Tunisie.  
E-mail: yousri.henchiri@univ-montp2.fr

ear processes in function spaces are developed in [1]. A comprehensive introduction to functional data analysis can be found in [25]. [12] focus on nonparametric approaches. Computational issues are explained in [26]. Nonparametric statistical methods for functional regression analysis, specifically the methods based on a Gaussian process prior in a functional space are discussed in [27]. [17] consider inferential procedures based on functional principal components. [31] mainly focuses on hypothesis testing problems about functional data. Among this, the functional linear model has received a special attention (see [24, 4, 5, 3, 15, 10, 6, 30] for main references).

In this paper, we are interested in the functional linear model

$$Y = \langle \theta, X \rangle + \varepsilon, \quad (1)$$

where  $\theta$  is the unknown function of the model,  $Y$  is a real variable of interest,  $\varepsilon$  is a centered real random variable representing the error of the model, with finite variance  $\mathbb{E}(\varepsilon^2) = \sigma_\varepsilon^2$ , and  $X$  is a functional covariate belonging to some functional space  $H$  endowed with an inner product  $\langle \cdot, \cdot \rangle$  and its associated norm  $\|\cdot\|$ . Usually,  $H$  is the space  $L^2([a, b])$  of square integrable functions defined on some real compact  $[a, b]$  and the corresponding inner product is defined by  $\langle f, g \rangle = \int_a^b f(t)g(t) dt$  for functions  $f, g \in L^2([a, b])$ . Without loss of generality, we consider our work on  $[0, 1]$ . Moreover, we assume that  $X$  and  $\varepsilon$  are independent.

All the previously cited works are devoted to analyze complete data, however, this is not the case in many interesting applications including for example survival data analysis. For this reason, we focus in this work on the problem of missing data (see [19, 14] for a wide introduction in the multivariate framework). This subject has been widely studied, in particular the way to impute missing data and the accuracy of this imputation according to the types of missing data: Missing Completely At Random (MCAR), Missing At Random (MAR) and Missing Not At Random (MNAR). Even if this problematic has received a lot of attention in a multivariate framework, it is not the case for the functional data framework. Our objective is to study the problem of combining regression imputation, missing data mechanisms and functional data analysis. As far as we know, few results are available for the moment. In MAR setting, [16] have explored this area by developing a functional multiple imputation approach modeling missing longitudinal response under a functional mixed effects model. They developed a Gibbs sampling algorithm to draw model parameters and imputations for missing values. Besides, [13] have considered two kinds of mean estimates of a scalar outcome, based on a sample in which an explanatory variable is observed for every subject while responses are missing (which is the closest to our context). A weak convergence result was proved. In MCAR setting, [23] have adapted a methodology based on the NIPALS (Nonlinear Iterative Partial Least Squares) algorithm, which provides

an imputation method for missing data, which have affected the functional covariates. In MNAR setting, [2] adapts a specification test for functional data with the presence of missing observations. His method is able to extract the information available in the observed portion of the data while being agnostic about the nature of the missing observations. In MAR and MCAR setting, [9] have recently proposed a nonparametric approach to missing value imputation and outlier detection for functional data. To our knowledge, there is no existing theoretical result in the case of functional linear model under missing assumption operating on the response variable, this problem only being until now the subject of studies in the multivariate framework (see for instance [20], [21]).

We carefully distinguish the missing data problem from a simple prediction problem. Indeed, the missing data mechanism involves a random variable (which indicates whether the response is missing or not) which plays a central role when obtaining our asymptotic results. This random variable and the variable  $X$  are dependent in the MAR case. This is also highlighted in [13]. In this paper, we first propose an imputation method, based on the completely observed cases, to replace missing values in the response of the functional linear model. We get mean square error rates for these imputed values. Secondly, once the database is completed, we are able to estimate the unknown function  $\theta$  of the model with the whole sample. This estimator can then be used for predicting other values of the response on a test set.

Combining missing data and functional variables offers a very large field of applications. Among all possible applications, environment is a core issue interesting many people for the future of our planet, in particular in the study of pollution indexes. The dataset we study here deals with temperature curves in some French cities to predict a specific pollution atmospheric index. The atmospheric index is missing in some cities in the northwest of France, for which the corresponding temperature curves (the explanatory variable) are mild, and leads to consider MAR data. The main objective is to get a map of the atmospheric index on the whole French territory.

The rest of the paper is organized as follows. Section 2 introduces the problem of functional linear model under missing assumption operating on the response variable and formulates our main results of the imputation method and of the mean square error for prediction of a new observation using the complete dataset. A simulation study is performed in Section 3. The environmental application is presented in Section 4. Some preliminary lemmas, which are used in the proofs of the main results, are collected in Section 5.

## 2 Imputation of a missing value of the response

### 2.1 Functional principal components regression

Let us consider a sample  $(X_i, Y_i)_{i=1, \dots, n}$  independent and identically distributed with the same distribution as  $(X, Y)$ . An estimation of  $\theta$  based on principal components analysis of the curves  $X_1, \dots, X_n$  has been studied in many papers, see for instance [4]. We recall below the construction of this estimator. Considering the covariance operator of  $X$  defined under the condition  $\mathbb{E}(\|X\|^2) < +\infty$  (which is supposed to be satisfied in the following) by

$$\Gamma u = \mathbb{E}(\langle X, u \rangle X),$$

for all  $u \in H$  and its empirical version

$$\widehat{\Gamma}_n u = \frac{1}{n} \sum_{i=1}^n \langle X_i, u \rangle X_i,$$

we call  $(\lambda_j)_{j \geq 1}$  (resp.  $(\widehat{\lambda}_j)_{j \geq 1}$ ) the sequence of eigenvalues of  $\Gamma$  (resp.  $\widehat{\Gamma}_n$ ) and  $(v_j)_{j \geq 1}$  (resp.  $(\widehat{v}_j)_{j \geq 1}$ ) the sequence of eigenfunctions of  $\Gamma$  (resp.  $\widehat{\Gamma}_n$ ). The identifiability of model (1) is ensured as long as we assume that  $\lambda_1 > \lambda_2 > \dots > 0$  (see [4]). Moreover, assuming that  $\widehat{\lambda}_{k_n} > 0$  for some integer  $k_n$  depending on  $n$ , the estimator of  $\theta$  is defined by

$$\widehat{\theta} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{k_n} \frac{\langle X_i, \widehat{v}_j \rangle Y_i}{\widehat{\lambda}_j} \widehat{v}_j. \quad (2)$$

A consistency result of this estimator is given in [4], while more recent results can be found in [3, 15]. In particular, [4] give technical conditions on the decreasing rate to zero of the eigenvalues  $\lambda_j$ 's in order to ensure the consistency of the estimator.

### 2.2 Operatorial point of view

We notice in this subsection that the model (1) can be seen from an operatorial point of view. Indeed, we can write the model

$$Y = \Theta X + \varepsilon, \quad (3)$$

where  $\Theta : H \rightarrow \mathbb{R}$  is a linear continuous operator defined by  $\Theta u = \langle \theta, u \rangle$  for any function  $u \in H$ . Let us consider  $\widehat{\Delta}_n$  the cross covariance operator defined by  $\widehat{\Delta}_n u = \frac{1}{n} \sum_{i=1}^n \langle X_i, u \rangle Y_i$ , for all  $u \in H$ . Then, it is easily seen that an estimator  $\widehat{\Theta}$  of  $\Theta$ , satisfying  $\widehat{\Theta} = \langle \widehat{\theta}, \cdot \rangle$ , is given by

$$\widehat{\Theta} = \widehat{\Pi}_{k_n} \widehat{\Delta}_n \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1}, \quad (4)$$

where  $\widehat{\Pi}_{k_n}$  is the projection operator onto the subspace  $\text{span}(\widehat{v}_1, \dots, \widehat{v}_{k_n})$ .

### 2.3 Imputation principle

Now, we present the context of missing data. There can be many reasons for which missing data can appear: breakdown in a measurement process, a person who is not willing to answer to some question of a questionnaire, ... We consider that some of the observations  $Y_1, \dots, Y_n$  are not available. We define the real variable  $\delta$  and we consider the sample  $(\delta_i)_{i=1, \dots, n}$  such that  $\delta_i = 1$  if the value  $Y_i$  is available and  $\delta_i = 0$  if the value  $Y_i$  is missing, for all  $i = 1, \dots, n$ . The data we observe are

$$\{(Y_i, \delta_i, X_i)\}_{i=1}^n.$$

We consider that the missing values are MAR. The MAR assumption implies that  $\delta$  and  $Y$  are conditionally independent given  $X$ . That is,

$$P(\delta = 1 \mid X, Y) = P(\delta = 1 \mid X). \quad (5)$$

Note that the MAR assumption is much weaker than MCAR (for which  $P(\delta = 1 \mid X, Y) = P(\delta = 1)$ ), as it allows the missing data to possibly depend on the observed data and may be reasonable for many practical problems. As a consequence of this MAR assumption, the variable  $\delta$  (the fact that an observation is missing) is independent of the error of the model  $\epsilon$ , conditionally on  $X$ . In the following, the number of missing values in the sample is denoted

$$m_n = \sum_{i=1}^n \mathbb{1}_{\{\delta_i=0\}}. \quad (6)$$

Then, to impute a missing value, say  $Y_\ell$  (where  $\ell$  is a given integer between 1 and  $n$ ), a simple way is to consider complete case analysis (see for instance [19, 7, 29, 22, 28]). This regression imputation method uses the pairs of observed data to define the estimator of the model coefficient. More precisely, we define

$$Y_{\ell, imp} = \frac{1}{n - m_n} \sum_{\substack{i=1 \\ i \neq \ell}}^n \sum_{j=1}^{k_n} \frac{\langle X_i, \widehat{v}_j \rangle \langle X_\ell, \widehat{v}_j \rangle \delta_i Y_i}{\widehat{\lambda}_j}. \quad (7)$$

From the operatorial point of view, the imputation of the missing value  $Y_\ell$  comes back to

$$Y_{\ell, imp} = \widehat{\Pi}_{k_n, obs} \widehat{\Delta}_{n, obs} \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell, \quad (8)$$

where  $\widehat{\Gamma}_{n,obs} = \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \cdot \rangle \delta_i X_i$ ,  $\widehat{\Delta}_{n,obs} = \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \cdot \rangle \delta_i Y_i$  and  $\widehat{\Pi}_{k_n,obs}$  is the projection operator onto the subspace  $\text{span}(\widehat{v}_{1,obs}, \dots, \widehat{v}_{k_n,obs})$  where  $\widehat{v}_{1,obs}, \dots, \widehat{v}_{k_n,obs}$  are the  $k_n$  first eigenfunctions of the covariance operator  $\widehat{\Gamma}_{n,obs}$ .

Now we give our main results. We consider the following assumptions.

(A.1) We assume that there exists a convex function  $\lambda$  such that  $\lambda(j) = \lambda_j$  for all  $j \geq 1$  that continuously interpolates the  $\lambda_j$ 's between  $j$  and  $j+1$ .

(A.2) There exists a positive constant  $C$  such that

$$\mathbb{E}(\|X\|^4) \leq C.$$

Our assumptions are quite classic in this context. Assumption (A.1) is similar to an assumption from [11]. It is a mild condition that allows a large class of decreasing rate of eigenvalues for the covariance operator  $\Gamma$ , for example polynomial decay or exponential decay (see example 1 below, in page 7, for more details). Assumption (A.2) holds for many processes  $X$  (Gaussian processes, bounded processes) and can also be found for example in [4]. Then, we give our main results.

*Remark 1* Notice that the assumptions (A.1) and (A.2) are just needed to obtain a convergence rate, whether there are missing data on the response or not. The only assumption needed on missing data is actually the MAR model.

**Theorem 1** *Assume (A.1) and (A.2) are satisfied, if, moreover  $\lambda_{k_n} k_n$  goes to zero as  $n$  goes to infinity, we have*

$$\mathbb{E}\left(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle\right)^2 = \sum_{j=k_n+1}^{+\infty} \left(\Theta \Gamma^{1/2} v_j\right)^2 + \frac{\sigma_{\varepsilon}^2 k_n}{n-m_n} + o\left(\frac{k_n}{n-m_n}\right).$$

Moreover, for the aggregate error of all the imputed values, we have

$$\sum_{\ell=1}^n (1-\delta_{\ell}) \mathbb{E}\left(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle\right)^2 = m_n \sum_{j=k_n+1}^{+\infty} \left(\Theta \Gamma^{1/2} v_j\right)^2 + \frac{\sigma_{\varepsilon}^2 k_n m_n}{n-m_n} + o\left(\frac{k_n m_n}{n-m_n}\right).$$

In order to precise the convergence rate of the imputed value  $Y_{\ell,imp}$  to the real one  $\langle \theta, X_{\ell} \rangle$ , we need an additional notation. For a function  $\varphi : \mathbb{R}_{+}^* \rightarrow \mathbb{R}_{+}^*$  and a positive real number  $L$ , we define

$$\mathcal{C}(\varphi, L) = \left\{ T : H \rightarrow \mathbb{R} \ / \ \forall j \geq 1, T v_j \leq L \sqrt{\varphi(j)} \right\}.$$

Note that simple cases satisfy the fact that  $\Theta \Gamma^{1/2}$  belongs to  $\mathcal{C}(\varphi, L)$ . For example, consider the operator  $\Theta$  expressed in the eigenfunctions basis  $(v_j)_{j \geq 1}$

such that  $\Theta u = \sum_{j=1}^{+\infty} \theta_j \langle v_j, u \rangle$  for any  $u \in H$ , with  $\theta_j$  going to zero as  $j$  goes to infinity. Hence there exists a bound  $L$  such that  $\theta_j \leq L$  for any  $j \geq 1$  and  $\Theta \Gamma^{1/2} v_j = \theta_j \sqrt{\lambda_j} \leq L \sqrt{\lambda_j}$ .

*Remark 2* We introduce two notations to compare the magnitudes of two functions  $\tilde{u}(x)$  and  $\tilde{v}(x)$  as the argument  $x$  tends to a limit  $\tilde{\ell}$  (not necessarily finite). The notation  $\tilde{u}(x) \underset{x \rightarrow \tilde{\ell}}{\sim} \tilde{v}(x)$ , stands for

$$\lim_{x \rightarrow \tilde{\ell}} \frac{\tilde{u}(x)}{\tilde{v}(x)} = 1,$$

and the notation  $\tilde{u}(x) \underset{x \rightarrow \tilde{\ell}}{\lesssim} \tilde{v}(x)$  denotes that  $|\tilde{u}(x)/\tilde{v}(x)|$  remains bounded as  $x \rightarrow \tilde{\ell}$ .

**Theorem 2** Let  $L = \|\Theta \Gamma^{1/2}\|_{\infty}$  and  $\varphi$  the function defined by  $\varphi(j) = \frac{(\Theta \Gamma^{1/2} v_j)^2}{L^2}$  for all  $j \geq 1$  that continuously interpolates the  $\varphi(j)$ 's between  $j$  and  $j+1$ . Under assumptions (A.1)-(A.2), the operator  $\Theta \Gamma^{1/2}$  belongs to  $\mathcal{C}(\varphi, L)$  and

$$\mathbb{E} (Y_{\ell, imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2\sigma_{\varepsilon}^2 \frac{k_n^*}{n - m_n},$$

where  $k_n^*$  is the solution of the equation in  $x$

$$\int_x^{+\infty} \varphi(t) dt = \frac{\sigma_{\varepsilon}^2}{L^2(n - m_n)} x. \quad (9)$$

Again, for the aggregate error of all the imputed values, we have

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E} (Y_{\ell, imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2\sigma_{\varepsilon}^2 \frac{k_n^* m_n}{n - m_n}.$$

*Remark 3* Notice that the equation (9) has a unique solution (the left and right hand sides are decreasing and increasing in  $x$ , respectively).

The last result giving the convergence rate of the imputed value  $Y_{\ell, imp}$  is similar to the convergence rate obtained in [11] (who considered the case of a completely observed functional response). The rate is simply affected by the number  $m_n$  of missing values. We precise the resulting rate of convergence in the following example.

*Example 1* We consider three different cases.



[Case 1] If the number of missing values is negligible compared to the sample size, that is  $m_n = a_n n$  with  $a_n$  going to zero as  $n$  goes to infinity, then the result of Theorem 2 becomes

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2\sigma_{\varepsilon}^2 \frac{k_n^*}{n},$$

for a single imputation and

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2\sigma_{\varepsilon}^2 \frac{k_n^* m_n}{n},$$

for the aggregate error of all the imputed values.

To go further, we consider two different functions  $\varphi$  such that  $\varphi_{\text{pol}}(j) = C_{\alpha} j^{-(2+\alpha)}$  and  $\varphi_{\text{exp}}(j) = D_{\alpha} \exp(-\alpha j)$  where  $C_{\alpha}$  and  $D_{\alpha}$  are positive constants and  $\alpha > 0$ . Then the solution of equation (9) is

$$\begin{cases} k_{n,pol}^* \underset{n \rightarrow +\infty}{\sim} \left( \frac{C_{\alpha} L^2}{(1+\alpha)\sigma_{\varepsilon}^2} \right)^{1/(2+\alpha)} n^{1/(2+\alpha)}, & \text{if } \varphi = \varphi_{\text{pol}}, \\ k_{n,exp}^* \underset{n \rightarrow +\infty}{\lesssim} \frac{\log n}{\alpha}, & \text{if } \varphi = \varphi_{\text{exp}}. \end{cases}$$

For  $\varphi = \varphi_{\text{pol}}$ , the result of Theorem 2 becomes

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2(\sigma_{\varepsilon}^2)^{(1+\alpha)/(2+\alpha)} \left( \frac{C_{\alpha} L^2}{1+\alpha} \right)^{1/(2+\alpha)} \frac{1}{n^{(1+\alpha)/(2+\alpha)}},$$

for a single imputation and

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2(\sigma_{\varepsilon}^2)^{(1+\alpha)/(2+\alpha)} \left( \frac{C_{\alpha} L^2}{1+\alpha} \right)^{1/(2+\alpha)} a_n n^{1/(2+\alpha)},$$

for the aggregate error of all the imputed values.

For  $\varphi = \varphi_{\text{exp}}$ , the result of Theorem 2 becomes

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\lesssim} \frac{\log n}{\alpha n},$$

for a single imputation and

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\lesssim} \frac{a_n \log n}{\alpha},$$

for the aggregate error of all the imputed values. In particular, we can see that missing data do not affect the convergence rate for a single imputed value. The rate  $\frac{1}{n^{(1+\alpha)/(2+\alpha)}}$  matches the usual optimal rates in this context.

The rate  $\frac{\log n}{\alpha n}$  is not exact but obviously sharp since parametric up to a logarithm. For the aggregate error of several imputed values, the number of missing values plays a crucial role, since the convergence depends on the fact that  $a_n n^{1/(2+\alpha)}$  or  $a_n \log n$  go to zero as  $n$  goes to infinity.

[Case 2] If the number of missing values is proportional to the sample size,  $m_n = \rho n$  with  $0 < \rho < 1$ , then the result of Theorem 2 becomes

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 \frac{\sigma_{\varepsilon}^2}{1 - \rho} \frac{k_n^*}{n},$$

for a single imputation and

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 \frac{\sigma_{\varepsilon}^2 \rho}{1 - \rho} k_n^*,$$

for the aggregate error of all the imputed values. Then, for a single imputation, the rate of convergence of Theorem 2 is

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 (\sigma_{\varepsilon}^2)^{(1+\alpha)/(2+\alpha)} \left( \frac{C_{\alpha}(1 - \rho)L^2}{1 + \alpha} \right)^{1/(2+\alpha)} \frac{1}{n^{(1+\alpha)/(2+\alpha)}},$$

for  $\varphi = \varphi_{\text{pol}}$  and

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\lesssim} \frac{\log n}{\alpha n},$$

for  $\varphi = \varphi_{\text{exp}}$ . In this case, missing data do not affect the convergence rate of a single imputed value. However, it affects the convergence of the aggregate error term for several imputed values, since the term of order  $k_n^*$  cannot converge to zero.

[Case 3] If the number of observed values is  $u_n = o(n)$ , then the result of Theorem 2 becomes

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 \sigma_{\varepsilon}^2 \frac{k_n^*}{u_n},$$

for a single imputation and

$$\sum_{\ell=1}^n (1 - \delta_{\ell}) \mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 \sigma_{\varepsilon}^2 \frac{k_n^* m_n}{u_n},$$

for the aggregate error of all the imputed values. Then, for a single imputation, the rate of convergence of Theorem 2 is

$$\mathbb{E}(Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\sim} 2 (\sigma_{\varepsilon}^2)^{(1+\alpha)/(2+\alpha)} \left( \frac{C_{\alpha} L^2}{1 + \alpha} \right)^{1/(2+\alpha)} \frac{1}{u_n^{(1+\alpha)/(2+\alpha)}},$$

for  $\varphi = \varphi_{\text{pol}}$  and

$$\mathbb{E} (Y_{\ell,imp} - \langle \theta, X_{\ell} \rangle)^2 \underset{n \rightarrow +\infty}{\lesssim} \frac{\log u_n}{\alpha u_n},$$

$\varphi = \varphi_{\text{exp}}$ . In this case, missing data affects the convergence rate of a single imputed value. This seems natural since the number of missing values is much more important than the number of observed values. The convergence results collapse for the aggregate error of several imputed values with a term of order  $\frac{k_n^* m_n}{u_n}$ .

#### 2.4 Estimation of $\theta$ and prediction of future values

Once the database being reconstructed, we can use the full database to estimate the functional coefficient  $\theta$  of the model (directly inspired from (2)) (see also [8]), namely

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{k_n} \frac{\langle X_i, \hat{v}_j \rangle Y_i^*}{\hat{\lambda}_j} \hat{v}_j, \quad (10)$$

where  $Y_i^* = Y_i \delta_i + Y_{i,imp}(1 - \delta_i)$  for all  $i = 1, \dots, n$ . Then this estimator of  $\theta$  can be used to predict new values of the response  $Y$  on a test sample. Indeed, if  $X_{new}$  is a new curve, the corresponding predicted response value is

$$\hat{Y}_{new} = \langle X_{new}, \tilde{\theta} \rangle = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{k_n} \frac{\langle X_i, \hat{v}_j \rangle \langle X_{new}, \hat{v}_j \rangle Y_i^*}{\hat{\lambda}_j}. \quad (11)$$

We give below a result allowing to control the mean square prediction error of  $\hat{Y}_{new}$ .

**Theorem 3** *Under the assumptions of Theorem 1, if we additionally assume that  $m_n = o(n)$  and  $m_n^2 k_n = O(n)$ , then*

$$\mathbb{E} \left( \hat{Y}_{new} - \langle \theta, X_{new} \rangle \right)^2 = \sum_{j=k_n+1}^{+\infty} \left( \Theta \Gamma^{1/2} v_j \right)^2 + O \left( \frac{k_n}{n} \right).$$

*Remark 4* This result shows that, under the condition that there are not too many missing values, the convergence rate of the mean square error prediction of a new value of the covariate remains the same compared to the non missing values case.

### 3 Simulations

To observe the behavior of our estimator in practice, this section considers a simulation study.

### 3.1 Models

Two models are considered:

$$Model_1 : Y = \int_0^1 \sin(4\pi t) X_t dt + \epsilon, \quad (12)$$

$$Model_2 : Y = \int_0^1 (\log(15t^2 + 10) + \cos(4\pi t)) X_t dt + \epsilon, \quad (13)$$

where the error  $\epsilon$  is a Gaussian noise :  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon)$  and

- in equation (12),  $X := \{X_t\}_{t \in [0,1]}$  is the standard Brownian motion.
- In equation (13),  $X := \{X_t\}_{t \in [0,1]}$  is a Gaussian process where the covariance function is defined as  $cov(X_t, X_{t'}) = \exp(-\frac{|t-t'|^2}{0.2})$ .

The simulation aims at considering processes  $X$  with different regularities (the standard Brownian motion being the case of the less smooth) in order to see if it has an impact on the results.

All the procedures described below were implemented by using the R software:

- ★ the trajectories of  $X_i$ ,  $1 \leq i \leq n$ , in the two models are discretized in  $p = 100$  equidistant points,
- ★ values of  $Y$  are computed using integration by rectangular interpolation,
- ★ the variability of noise is such that  $\sigma_\epsilon = \tau * \text{Var}\left(\int_0^1 \theta(t) X(t) dt\right) \approx 0.2$ ,
- ★ the sample sizes are respectively  $n = 100, 300$  and  $1200$  for the training sets  $(X_1, Y_1), \dots, (X_n, Y_n)$  and  $n_1 = 50, 150$  and  $600$  for the test sets  $(X_{n+1}, Y_{n+1}), \dots, (X_{n+n_1}, Y_{n+n_1})$ .

Note that some Monte Carlo experiments are achieved to determine the values of  $\tau$ :  $\tau \approx 21.726$  for the *model*<sub>1</sub> and  $\tau \approx 0.048$  for the *model*<sub>2</sub>.

### 3.2 Criteria

The criteria we used are the following. Criteria 1, 2, 3 are related to the imputation step with the training samples, criteria 4, 5, 6 are related to the prediction step with the test samples, and criteria 7 is related to the estimation step with the reconstructed database.

- Criterion 1: the mean square errors (*MSE*) averaged over  $\mathcal{S}$  samples

$$\overline{MSE} = \frac{1}{\mathcal{S}} \sum_{j=1}^{\mathcal{S}} MSE(j),$$

where  $MSE(j) = \frac{1}{m_n} \sum_{\ell=1}^n (Y_{\ell, imp}^j - \langle \theta, X_{\ell}^j \rangle)^2 (1 - \delta_{\ell})$  is the mean square error computed on the  $j^{th}$  simulated sample,  $j \in \{1, \dots, \mathcal{S}\}$ .

- Criterion 2: the mean absolute errors ( $MAE$ ) averaged over  $\mathbf{S}$  samples

$$\overline{MAE} = \frac{1}{\mathbf{S}} \sum_{j=1}^{\mathbf{S}} MAE(j),$$

where  $MAE(j) = \frac{1}{m_n} \sum_{\ell=1}^n |Y_{\ell,imp}^j - \langle \theta, X_{\ell}^j \rangle| (1 - \delta_{\ell})$  is the mean absolute error computed on the  $j^{th}$  simulated sample.

- Criterion 3: the ratio between the mean square prediction error and the mean square prediction error when the true mean is known averaged over  $\mathbf{S}$  samples

$$\overline{CR3} = \frac{1}{\mathbf{S}} \sum_{j=1}^{\mathbf{S}} CR3(j),$$

where  $CR3(j) = \frac{\sum_{\ell=1}^n (Y_{\ell,imp}^j - \langle \theta, X_{\ell}^j \rangle)^2 (1 - \delta_{\ell})}{\sum_{\ell=1}^n (\epsilon_{\ell}^j)^2 (1 - \delta_{\ell})}$  is the ratio between the

mean square prediction error and the mean square prediction error when the true mean is known, computed on the  $j^{th}$  simulated sample.

- Criterion 4: the mean square errors ( $MSE'$ ) averaged over  $\mathbf{S}$  samples

$$\overline{MSE'} = \frac{1}{\mathbf{S}} \sum_{j=1}^{\mathbf{S}} MSE'(j),$$

where  $MSE'(j) = \frac{1}{n_1} \sum_{\ell'=n+1}^{n+n_1} (Y_{\ell'}^j - \langle \theta, X_{\ell'}^j \rangle)^2$  is the mean square error computed on the  $j^{th}$  simulated sample,  $j \in \{1, \dots, \mathbf{S}\}$ .

- Criterion 5: the mean absolute errors ( $MAE'$ ) averaged over  $\mathbf{S}$  samples

$$\overline{MAE'} = \frac{1}{\mathbf{S}} \sum_{j=1}^{\mathbf{S}} MAE'(j),$$

where  $MAE'(j) = \frac{1}{n_1} \sum_{\ell'=n+1}^{n+n_1} |Y_{\ell'}^j - \langle \theta, X_{\ell'}^j \rangle|$  is the mean absolute error computed on the  $j^{th}$  simulated sample.

- Criterion 6: the ratio between the mean square prediction error and the mean square prediction error when the true mean is known averaged over  $\mathbf{S}$  samples

$$\overline{CR3'} = \frac{1}{\mathbf{S}} \sum_{j=1}^{\mathbf{S}} CR3'(j),$$

where  $CR3'(j) = \frac{\sum_{\ell'=n+1}^{n+n_1} (Y_{\ell'}^j - \langle \theta, X_{\ell'}^j \rangle)^2}{\sum_{\ell'=n+1}^{n+n_1} (\epsilon_{\ell'}^j)^2}$  is the ratio between the mean

square prediction error and the mean square prediction error when the true mean is known, computed on the  $j^{\text{th}}$  simulated sample.

- Criterion 7: the mean square errors ( $MSE''$ ) averaged over  $\mathcal{S}$  samples

$$\overline{MSE''} = \frac{1}{\mathcal{S}} \sum_{j=1}^{\mathcal{S}} MSE''(j),$$

where  $MSE''(j) = \|\tilde{\theta}^j - \theta\|^2$  is the square error of estimation computed on the  $j^{\text{th}}$  simulated sample. The  $MSE''$  criterion is decomposed into variance and square bias in our results.

Notice that all the criteria tend to zero when the sample size tends to infinity.

### 3.3 Methodology

We smooth the estimator (2) by a pre-processing step based on the Smooth Principal Components Regression (SPCR) [5]. In our context, we use the regression spline such as the original curves  $X_1, \dots, X_n$  are projected on a regression spline basis. Then, our estimator depends on other additional parameters: the number ' $\kappa$ ' of knots of the spline functions, the degree ' $q$ ' of spline functions and the number ' $m$ ' of derivatives. Here, we have fixed the number of knots to be 20, the degree has been chosen to be 3 and the number of derivatives was fixed to the moderate value of 2. These parameters are not the most important in our study, especially in comparison with the choice of the number of principal components.

In this subsection, we show firstly how to determine the number of missing data in the MAR case. Secondly, we present a procedure to choose the optimal tuning parameter (the best dimension  $k_n^*$  of the projection space for the SPCR). Thirdly, in order to illustrate the performance of the estimator  $\tilde{\theta}$  using imputed values with the optimal chosen dimension, we have chosen a percentage of missing values equal to 45.8518% for  $model_1$  and equal to 46.8888% for  $model_2$  (we obtain this rate with  $ct = 1$  for both models, see next paragraph below).

### 3.3.1 Missing data simulation scenarios

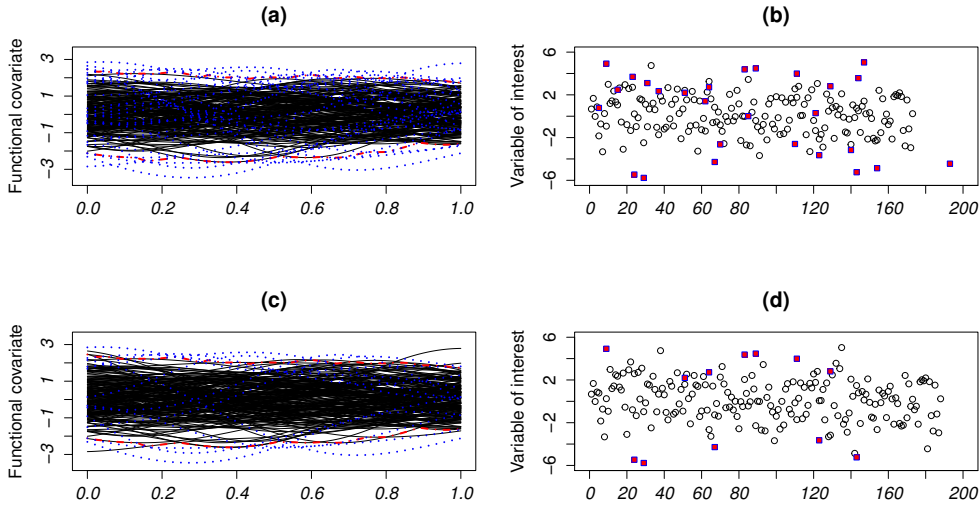
To determine the number of missing data in our simulations, we have adopted two scenarios. In the first one, we use an scheme based on the the confidence band associated to the functional covariate. This scheme is done in two ways (the first way illustrates the principle to determine missing data, the second way gives the possibility to control the missing data percentage). The second scenario is based on the logistic functional regression. We give below more details on both scenarios. In the simulations we present, we adopted the second scenario.

The fact that  $Y$  is observed or missing can be linked to a condition on the curve  $X$  which reaches high or low levels or not. In the first way, this number is associated to the number of curves which do not belong to some confidence band (90%, 95%, 97% and 99%). More precisely, each curve  $X_i$ ,  $i = 1, \dots, n$ , is said to belong to the confidence band if all discretization points are in the band. Then the variable of interest  $Y_i$ , associated to  $X_i$ , is called *available*. In the second way, we modify the first way such that a curve is said to belong to the confidence band if some rate of the discretization points (80 percent) are in the band. This strategy allows to control the rate of missing data, this rate being decreasing from the first way to the second. Fig. 1 illustrates a simple example. We have considered  $Model_2$  under  $n = 200$  observations and  $p = 100$  discretization points. For the first way, the number of missing data (27 points, see **(b)**) is associated to the number of curves that do not belong to the 97% confidence band (see **(a)**). In the second way, the number of missing data (12 points, see **(d)**) is associated to the number of curves that do not belong to the 97% confidence band (see **(c)**). Notice that, in the last case, if more than 20 discretization points of a curve  $X_i$  are not in the confidence band then the variable of interest  $Y_i$ , associated to  $X_i$ , is called *missing*.

Now, we present the second scenario which is a simpler strategy to simulate missing data. We used this simulation method to obtain the results we present. In the MAR case, we simulate  $\delta$  according to the logistic functional regression. The variable  $\delta$  follows the Bernoulli law with parameter  $p(X)$  such that

$$\log \left( \frac{p(X)}{1 - p(X)} \right) = \langle \alpha_0, X \rangle + ct,$$

where  $\alpha_0(t) = \sin(2\pi t)$  for all  $t \in [0, 1]$  and  $ct$  is a constant allowing to take different levels of missing data. We take  $ct = 2$  for around 12.5% of missing data,  $ct = 1$  for around 27.4% of missing data and  $ct = 0.2$  for around 44.9% of missing data. Notice that, in the MCAR case, we simulate  $\delta$  with the Bernoulli law with parameter  $p = 0.9$  (10% of missing data),  $p = 0.75$  (25% of missing data) or  $p = 0.6$  (40% of missing data).



**Fig. 1** Plots of functional covariate and variable of interest in the first way of scenario 1 (resp. (a) and (b)) and the second way of scenario 1 (resp. (c) and (d)) under Missing At Random case.

### 3.3.2 Criteria for optimal parameter selection

We focus on the procedure allowing to select the optimal tuning parameter. We consider a Generalized Cross Validation (GCV) criterion versus a Cross Validation (CV) criterion and K-fold Cross Validation (K-fold CV) criterion and we select the optimal tuning parameter  $k_n^*$  by minimizing these criteria. The GCV procedure is known to be computationally fast. The CV, K-fold CV and GCV criteria are respectively given as follows for imputation

$$\begin{aligned}
 CV(k_n) &= \frac{1}{n - m_n} \sum_{i=1}^n (\hat{Y}_i^{[-i]} - \langle \theta, X_i \rangle)^2 \delta_i, \\
 \text{K-fold CV}(k_n) &= \frac{1}{K} \sum_{k=1}^K |B_k|^{-1} \sum_{i \in B_k} (\hat{Y}_i^{[-B_k]} - \langle \theta, X_i \rangle)^2 \delta_i, \\
 \text{GCV}(k_n) &= \frac{(n - m_n) \sum_{i=1}^n (\hat{Y}_i - \langle \theta, X_i \rangle)^2 \delta_i}{((n - m_n) - k_n)^2}.
 \end{aligned}$$



The analogous criteria are given as follows for prediction

$$\begin{aligned} \text{CV}(k_n) &= \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i^{*[-i]} - \langle \theta, X_i \rangle)^2, \\ \text{K-fold CV}(k_n) &= \frac{1}{K} \sum_{k=1}^K |B_k|^{-1} \sum_{i \in B_k} (\hat{Y}_i^{*[-B_k]} - \langle \theta, X_i \rangle)^2, \\ \text{GCV}(k_n) &= \frac{n \sum_{i=1}^n (\hat{Y}_i^* - \langle \theta, X_i \rangle)^2}{(n - k_n)^2}, \end{aligned}$$

where  $\hat{Y}_i^{*[-i]}$  and  $\hat{Y}_i^{*[-B_k]}$  respectively mean that the value of  $Y_i$  is predicted using the whole sample except the  $i^{\text{th}}$  observation or except the set of observations indexed in  $B_k$ . In the same way  $\hat{Y}_i^{*[-i]}$  and  $\hat{Y}_i^{*[-B_k]}$  respectively mean that the value of  $Y_i$  is predicted using the whole sample except the  $i^{\text{th}}$  observation or except the set of observations indexed in  $B_k$ . The data set is randomly partitioned into  $K$  equally sized (as equal as possible) subsets  $\cup_{k=1}^K B_k$  such that  $B_j \cap B_k = \emptyset$  ( $j \neq k$ ). In practice, often  $K = 5$  or  $K = 10$  are used. In our case, the K-fold CV splits are chosen in a special deterministic way. For imputation, we consider

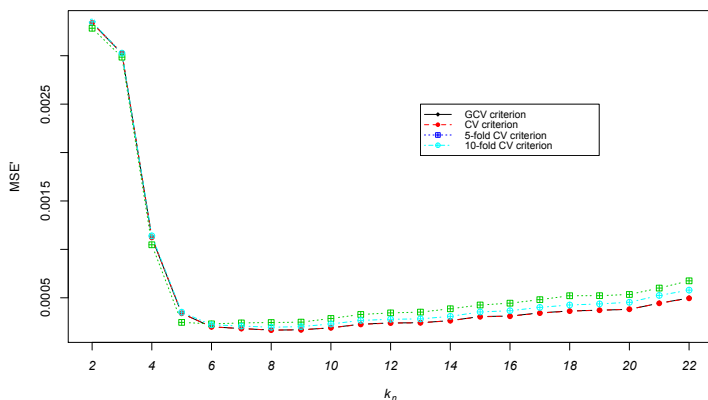
$$\text{K-fold CV}(k_n) = \frac{1}{K} \sum_{k=1}^K ((n - m_n)/K)^{-1} \sum_{i=(n(k-1))/K+1}^{nk/K} (\hat{Y}_i^{*[-k]} - \langle \theta, X_i \rangle)^2 \delta_i.$$

The analogous criterion is given as follows for prediction

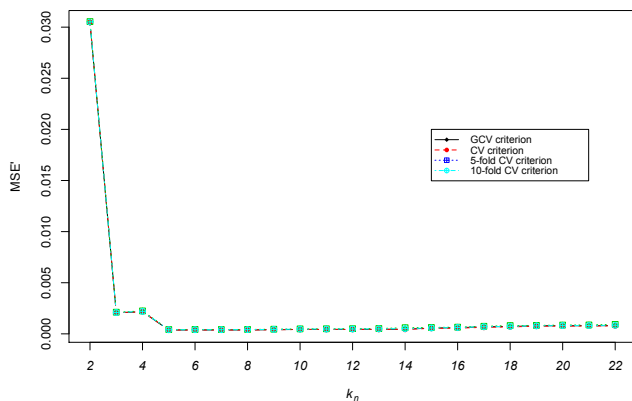
$$\text{K-fold CV}(k_n) = \frac{1}{K} \sum_{k=1}^K (n/K)^{-1} \sum_{i=(n(k-1))/K+1}^{nk/K} (\hat{Y}_i^{*[-k]} - \langle \theta, X_i \rangle)^2.$$

In order to illustrate the advantage of the GCV criterion, we compared the computational times to obtain the tuning parameter with the three criteria on a growing sequence of dimension  $k_n = 2, \dots, 22$ . The characteristics of the computer used to perform these computations were MacBook pro: Processor 2.66 GHz intel core 2 Duo, Memory 4 Gb 1067 MHz DDR3. The computational times are displayed in Table 15 in the appendix. The GCV criterion shows a clear advantage with regard to computational time compared with the CV and K-fold criteria. In addition, we see that the three criteria behave in the same way and select the same optimal projection dimension (see Fig. 2 and 3) for both models (under  $n = 1000$  and  $p = 100$ ).

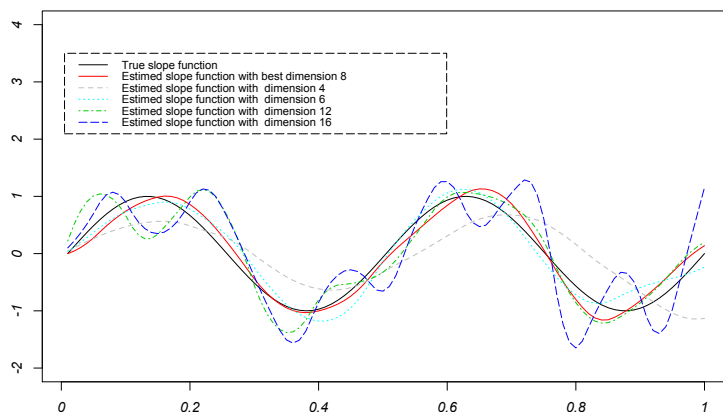
We show on Fig. 4 and Fig. 5 different estimates of the slope function of the *Model*<sub>1</sub> and *Model*<sub>2</sub> (under  $n = 1000$  and  $p = 100$ ) with different values of dimension ( $k_n = 4, 6, 8, 12, 16$ ) and ( $k_n = 2, 3, 5, 7, 8$ ), respectively, by using the GCV criterion (used for its computational efficiency).



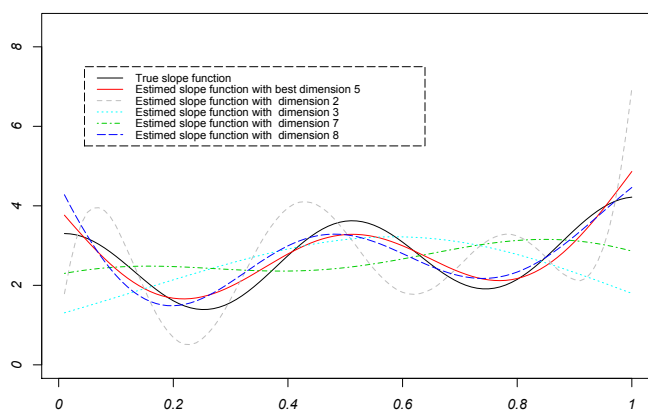
**Fig. 2** GCV, CV and K-fold criteria for different values of dimension  $k_n$  in  $model_1$ : best dimension  $k_n^* = 8$  and  $MSE' (\times 10^4) = 1.6640$  (in GCV criterion case), best dimension  $k_n^* = 6$  and  $MSE' (\times 10^4) = 2.3081$  (in 5-fold CV criterion case), best dimension  $k_n^* = 8$  and  $MSE' (\times 10^4) = 1.9584$  (in 10-fold CV criterion case), best dimension  $k_n^* = 8$  and  $MSE' (\times 10^4) = 1.6598$  (in CV criterion case), for the  $model_1$ .



**Fig. 3** GCV, CV and K-fold criteria for different values of dimension  $k_n$  in  $model_2$ : best dimension  $k_n^* = 5$  and  $MSE' (\times 10^4) = 3.7589$  (in GCV criterion case), best dimension  $k_n^* = 5$  and  $MSE' (\times 10^4) = 4.2132$  (in 5-fold CV criterion case), best dimension  $k_n^* = 5$  and  $MSE' (\times 10^4) = 3.9758$  (in 10-fold CV criterion case), best dimension  $k_n^* = 5$  and  $MSE' (\times 10^4) = 3.7270$  (in CV criterion case), for the  $model_2$ .



**Fig. 4** Plots of the true slope function (solid black) and estimates with different values of dimension  $k_n$  in *model*<sub>1</sub>. The plots of estimates slope function with best dimension  $k_n^* = 8$  (solid red), with dimension  $k_n = 4$  (dotted), with dimension  $k_n = 6$  (dashed), with dimension  $k_n = 12$  (dotdashed), with dimension  $k_n = 16$  (twodash).



**Fig. 5** Plots of the true slope function (solid black) and estimates with different values of dimension  $k_n$  in *model*<sub>2</sub>. The plots of estimates slope function with best dimension  $k_n^* = 5$  (solid red), with dimension  $k_n = 2$  (dotted), with dimension  $k_n = 3$  (dashed), with dimension  $k_n = 7$  (dotdashed), with dimension  $k_n = 8$  (twodash).

### 3.4 Analysis of results

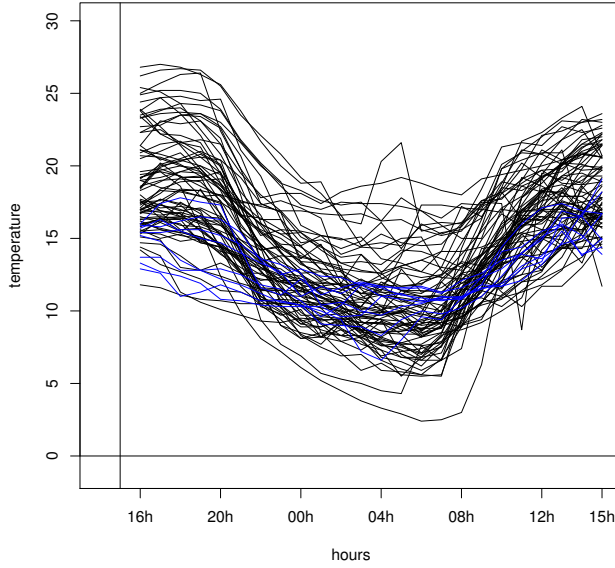
In this subsection, we analyse the results of the criteria presented in the previous subsection. Both MAR and MCAR context were considered. The different

results given in Appendix A. Tables 2, 3, 4, 5 give the mean and standard deviation errors for the imputed values on training samples for both models and for MAR and MCAR cases. Tables 6, 7, 8, 9 give the mean and standard deviation errors for the predicted values on test samples for both models and for MAR and MCAR cases. Tables 10, 11, 12, 13 give the mean and standard deviation errors for the estimation of  $\theta$  using the fulfilled database with imputed values for both models and for MAR and MCAR cases. We can see that the errors increase when the rate of missing data increases. Similarly, the errors decrease as the size of the sample increases. When we compare the case of MAR and MCAR, we see that the error in case of MAR is higher than in the MCAR case. Moreover, we can see that the regularity of the process  $X$  does not have a crucial impact on the results. All the results in these simulations are in accordance with what we can expect and confirm the theoretical results obtained in the previous section.

## 4 Application

In order to illustrate the contribution of our approach in functional prediction setting when the covariates are functions and some observations of the real response are missing, we present in this section an environmental dataset study.

We start by describing the dataset. The functional covariate  $X$  is a daily temperature curve in some cities in France (from May 7, 2015 at 4 pm up to May 8, 2015 at 3 pm) obtained from [www.meteociel.fr](http://www.meteociel.fr). This daily continuous curve is observed at some discretization points (here, at 24 discretization points, every hour). The graphical display of this daily temperature curves can be observed in Fig 6. The response variable  $Y$  is an atmospheric index of air quality called ATMO (for a detailed description of this atmospheric index, see [www.atmo-france.org](http://www.atmo-france.org)). Its values range from 1 (very good quality of air) to 10 (very bad quality of air). We obtained the values of the atmospheric index on May 8, 2015, for these same cities, from [www2.prevoir.org](http://www2.prevoir.org). Furthermore, we added some cities for which the temperature curve is available but the atmospheric index is missing. Notice that the response is missing for mild temperature curves cities: the fact that the value of the response variable  $Y$  is missing for these cities depends on the temperature curve  $X$ , and thus we consider the MAR case. We also refer the reader to the paper [18] for more discussions about missing data mechanism when dealing with air quality data. In particular, this paper highlights the fact that air quality missing data can be considered as MAR. Fig 7 illustrates the selected cities in our study, the blue cities are given when the response variable  $Y$  is missing and the red cities are given when the response variable  $Y$  is observed. It is of primary importance to get a map of the atmospheric index on the whole French territory, and thus to impute missing data.



**Fig. 6** Plot of the 78 daily temperature curves (the blue curves are given when the response variable  $Y$  is missing).

We have built a sample of 78 pairs  $\{(Y_i, X_i)\}_{i=1}^{78}$ , where we have 8 missing values of the variable  $Y$  (the  $Y_i$ 's,  $i = 71, \dots, 78$ , are missing). Our goal is to impute these missing values  $\{Y_i\}_{i=71}^{78}$ .

We have fixed the number of knots to be 20, the degree of splines has been chosen equal to 3 and the number of derivatives was fixed to the moderate value of 2. Then, we use the GCV criterion to find the best parameter of projection dimension  $k_n$  trying growing sequences:  $k_n = 2, 3, \dots, 21, 22$ . In order to see the impact of missing data on this dataset, we have randomly drawn 700 tests samples in the initial sample and computed prediction errors on these tests samples, using the remaining of the sample as training sample. Results are given in Table 14. Here again, the more we have missing data in the training set, the more the prediction error on the test sample is.

Now, we come back to the initial goal, imputing the missing data. The minimum value of the GCV criterion is reached for  $k_n^* = 5$  and  $MSE' (\times 10^2) = 20.791$ . Table 1 gives the imputed values of the missing data. We see imputed values mainly around 4, which is a moderate value for the atmospheric index corresponding to a good quality of air. It is in accordance with the fact that these cities have moderate temperature curves. We can mention two particular cases. The highest imputed value (4.161) corresponds to the city of Angers,

ml - Google Maps

<https://www.google.fr/maps/@47.4202487,3.5858941,6z/data=!4m2!6m1!1szHYkmp8S1XA.kzQ2zr...>

2 sur 3

20/10/2015 13:16

**Fig. 7** Map of France locating the selected cities of our study: the cities are red when the variable  $Y$  is observed and the cities are blue when the variable  $Y$  is missing.

and in parallel, we can see that the temperature curve of this city becomes high at the end of May 8. On the contrary, the lowest imputed value (3.491) corresponds to the city of Quimper, and the temperature curve of this city presents few variations along the 24 hours.

**Table 1** Imputed values of the missing response variable.

Missing values of $Y$	$Y_{71}$	$Y_{72}$	$Y_{73}$	$Y_{74}$	$Y_{75}$	$Y_{76}$	$Y_{77}$	$Y_{78}$
Imputed values	4.161	3.496	3.850	3.758	3.590	3.491	3.990	3.821

## 5 Proof of the results

### 5.1 Proof of Theorem 1

We begin with the following decomposition

$$\hat{\Delta}_{n,obs} = \frac{1}{n - m_n} \sum_{i=1}^n \langle X_i, \cdot \rangle \delta_i \theta X_i + \frac{1}{n - m_n} \sum_{i=1}^n \langle X_i, \cdot \rangle \delta_i \varepsilon_i = \theta \hat{\Gamma}_{n,obs} + U_{n,obs},$$

with  $U_{n,obs} = \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \cdot \rangle \delta_i \varepsilon_i$ . Then,  $\varepsilon$  being independent from  $X$  and  $\delta$  (MAR assumption), we deduce

$$\begin{aligned} \mathbb{E} (Y_{\ell,imp} - \langle \theta, X_\ell \rangle)^2 &= \mathbb{E} \left( \Theta \widehat{\Pi}_{k_n,obs} X_\ell - \Theta X_\ell \right)^2 \\ &\quad + \mathbb{E} \left( \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n,obs} \widehat{\Gamma}_{n,obs} \right)^{-1} X_\ell \rangle \delta_i \varepsilon_i \right)^2 \\ &\leq 2\mathbb{E} \left( \Theta \widehat{\Pi}_{k_n,obs} X_\ell - \Theta \Pi_{k_n,obs} X_\ell \right)^2 \\ &\quad + 2\mathbb{E} \left( \Theta \Pi_{k_n,obs} X_\ell - \Theta X_\ell \right)^2 \\ &\quad + \mathbb{E} \left( \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n,obs} \widehat{\Gamma}_{n,obs} \right)^{-1} X_\ell \rangle \delta_i \varepsilon_i \right)^2, \end{aligned}$$

where  $\Pi_{k_n,obs}$  is the projection onto the subspace  $\text{span}(v_{1,obs}, \dots, v_{k_n,obs})$  where  $v_{1,obs}, \dots, v_{k_n,obs}$  are the  $k_n$  first eigenfunctions of the covariance operator  $\Gamma_{n,obs}$ . For a single imputation, the end of the proof of Theorem 1 is based on the following lemmas. For the aggregate error term of  $m_n$  imputed values, it is just a sum of  $m_n$  terms that behave like the term for single imputation.

**Lemma 1** *We have*

$$\mathbb{E} \left( \Theta \widehat{\Pi}_{k_n,obs} X_\ell - \Theta \Pi_{k_n,obs} X_\ell \right)^2 = o \left( \frac{\lambda_{k_n} k_n^2}{n-m_n} + \frac{k_n}{n-m_n} \right).$$

**Lemma 2** *We have*

$$\mathbb{E} \left( \frac{1}{n-m_n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n,obs} \widehat{\Gamma}_{n,obs} \right)^{-1} X_\ell \rangle \delta_i \varepsilon_i \right)^2 = \frac{\sigma_\varepsilon^2 k_n}{n-m_n} + o \left( \frac{k_n}{n-m_n} \right).$$

**Lemma 3** *We have*

$$\mathbb{E} \left( \Theta \Pi_{k_n,obs} X_\ell - \Theta X_\ell \right)^2 = \sum_{j=k_n+1}^{+\infty} \left( \Theta \Gamma^{1/2} v_j \right)^2.$$

## 5.2 Proof of Lemma 1

Writing  $X_\ell$  in the basis  $(v_j)_{j \geq 1}$ , we obtain

$$\begin{aligned} &\mathbb{E} \left( \Theta \widehat{\Pi}_{k_n,obs} X_\ell - \Theta \Pi_{k_n,obs} X_\ell \right)^2 \\ &= \sum_{j=1}^{+\infty} \sum_{j'=1}^{+\infty} \mathbb{E} \left[ \langle X_\ell, v_j \rangle \langle X_\ell, v_{j'} \rangle \Theta \left( \widehat{\Pi}_{k_n,obs} - \Pi_{k_n,obs} \right) v_j \Theta \left( \widehat{\Pi}_{k_n,obs} - \Pi_{k_n,obs} \right) v_{j'} \right]. \end{aligned}$$

Noticing that the variable  $X_\ell$  corresponds to the missing data  $Y_\ell$  hence independent of  $\hat{\Pi}_{k_n, obs}$ , we get

$$\begin{aligned} & \mathbb{E} \left( \Theta \hat{\Pi}_{k_n, obs} X_\ell - \Theta \Pi_{k_n, obs} X_\ell \right)^2 \\ &= \sum_{j=1}^{+\infty} \sum_{j'=1}^{+\infty} \langle \Gamma v_j, v_{j'} \rangle \mathbb{E} \left[ \Theta \left( \hat{\Pi}_{k_n, obs} - \Pi_{k_n, obs} \right) v_j \Theta \left( \hat{\Pi}_{k_n, obs} - \Pi_{k_n, obs} \right) v_{j'} \right] \\ &= \sum_{j=1}^{+\infty} \lambda_j \mathbb{E} \left[ \Theta \left( \hat{\Pi}_{k_n, obs} - \Pi_{k_n, obs} \right) v_j \right]^2. \end{aligned}$$

Now, following the proof of Proposition 15 in [11], for any  $m \geq 1$  we denote  $\mathcal{B}_m$  the oriented circle of the complex plane with center  $\lambda_m$  and radius  $\rho_m/2$  where  $\rho_m = \min(\lambda_m - \lambda_{m+1}, \lambda_{m-1} - \lambda_m)$  for  $m \geq 2$  and  $\rho_1 = \lambda_2 - \lambda_1$ . With the convexity assumption (A.1), we actually have  $\rho_m = \lambda_m - \lambda_{m+1}$  for all  $m \geq 1$ . With these notations, denoting by  $\iota$  the complex number such that  $\iota^2 = -1$ , the difference between the projection operators  $\hat{\Pi}_{k_n, obs}$  and  $\Pi_{k_n, obs}$  can be written

$$\hat{\Pi}_{k_n, obs} - \Pi_{k_n, obs} = \frac{1}{2\pi\iota} \sum_{m=1}^{k_n} \int_{\mathcal{B}_m} \Lambda(z) \left( \Gamma - \hat{\Gamma}_{n, obs} \right) \Lambda(z) dz,$$

where  $\Lambda(z) = (zI - \Gamma)^{-1}$ . Noticing that  $\Lambda(z)v_j = \frac{1}{z - \lambda_j} v_j$ , we deduce

$$\begin{aligned} & \Theta \left( \hat{\Pi}_{k_n, obs} - \Pi_{k_n, obs} \right) v_j \\ &= \frac{1}{2\pi\iota} \sum_{m=1}^{k_n} \Theta \int_{\mathcal{B}_m} \Lambda(z) \left( \Gamma - \hat{\Gamma}_{n, obs} \right) \frac{dz}{z - \lambda_j} \\ &= \frac{1}{2\pi\iota} \sum_{m=1}^{k_n} \Theta \int_{\mathcal{B}_m} \sum_{j'=1}^{+\infty} \frac{\langle \left( \Gamma - \hat{\Gamma}_{n, obs} \right) v_j, v_{j'} \rangle v_{j'}}{(z - \lambda_{j'})(z - \lambda_j)} dz. \end{aligned}$$

Still using the results from [11], we have

$$\sum_{m=1}^{k_n} \int_{\mathcal{B}_m} \frac{dz}{(z - \lambda_{j'})(z - \lambda_j)} = \begin{cases} 0, & \text{if } j, j' > k_n, \\ 0, & \text{if } j, j' \leq k_n, \\ (\lambda_j - \lambda_{j'})^{-1}, & \text{if } j \leq k_n < j', \\ (\lambda_{j'} - \lambda_j)^{-1}, & \text{if } j' \leq k_n < j. \end{cases}$$

hence we deduce



$$\begin{aligned}
& \mathbb{E} \left( \Theta \widehat{\Pi}_{k_n, obs} X_\ell - \Theta \Pi_{k_n, obs} X_\ell \right)^2 \\
&= \mathbb{E} \left[ \frac{1}{4\pi^2} \sum_{j=1}^{k_n} \lambda_j \left( \sum_{j'=k_n+1}^{+\infty} \frac{\langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_{j'} \rangle}{\lambda_j - \lambda_{j'}} \Theta v_{j'} \right)^2 \right] \\
&+ \mathbb{E} \left[ \frac{1}{4\pi^2} \sum_{j=k_n+1}^{+\infty} \lambda_j \left( \sum_{j'=1}^{k_n} \frac{\langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_{j'} \rangle}{\lambda_{j'} - \lambda_j} \Theta v_{j'} \right)^2 \right].
\end{aligned}$$

In the following,  $C$  corresponds to a generic constant. We denote  $\mathbb{E}(A)$  and  $\mathbb{E}(B)$  the above two terms. We start with the computation of  $\mathbb{E}(A)$ . Using the same technique as in [11], we get the following bound

$$\mathbb{E} \left( \langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_{j'} \rangle \langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_r \rangle \right) \leq \frac{C}{n - m_n} \lambda_j \sqrt{\lambda_{j'}} \sqrt{\lambda_r},$$

noticing that the  $n$  rate of convergence given in [11] is here transformed into the  $n - m_n$  rate because we use  $\widehat{\Gamma}_{n, obs}$  with  $n - m_n$  observed data. Hence we deduce

$$\begin{aligned}
& \mathbb{E} \left( \frac{\langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_{j'} \rangle}{\lambda_j - \lambda_{j'}} \Theta v_{j'} \right)^2 \\
&= \sum_{j'=k_n+1}^{+\infty} \sum_{r=k_n+1}^{+\infty} \frac{\mathbb{E} \left( \langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_{j'} \rangle \langle (\Gamma - \widehat{\Gamma}_{n, obs}) v_j, v_r \rangle \right)}{(\lambda_j - \lambda_{j'}) (\lambda_j - \lambda_r)} \Theta v_{j'} \Theta v_r \\
&\leq \frac{C \lambda_j}{n - m_n} \left( \sum_{j'=k_n+1}^{+\infty} \frac{\sqrt{\lambda_j}}{\lambda_j - \lambda_{j'}} \Theta v_{j'} \right)^2.
\end{aligned}$$

Coming back to the computation of  $\mathbb{E}(A)$ , we can write (using Lemma 12 in [11])

$$\begin{aligned}
\mathbb{E}(A) &\leq \frac{C}{n - m_n} \sum_{j=1}^{k_n} \frac{\lambda_j^2 \lambda_{k_n+1}}{(\lambda_j - \lambda_{k_n+1})^2} \left( \sum_{j'=k_n+1}^{+\infty} \Theta v_{j'} \right)^2 \\
&\leq \frac{C \lambda_{k_n+1}}{n - m_n} \sum_{j=1}^{k_n} \frac{(k_n + 1)^2}{(k_n + 1 - j)^2} \left( \sum_{j'=k_n+1}^{+\infty} \Theta v_{j'} \right)^2 \\
&\leq \frac{C \lambda_{k_n+1} (k_n + 1)^2}{n - m_n} \sum_{j=1}^{k_n} \frac{1}{j^2} \left( \sum_{j'=k_n+1}^{+\infty} \Theta v_{j'} \right)^2.
\end{aligned}$$

As  $\theta \in L^2([0, 1])$  (hence  $\theta$  is integrable), we finally get

$$\mathbb{E}(A) \leq \frac{C \lambda_{k_n} k_n^2}{n - m_n} a_n,$$

where  $(a_n)_{n \geq 1}$  is a sequence of real numbers going to zero as  $n$  goes to infinity. We are now interested in the computation of  $\mathbb{E}(B)$ . Beginning in the same way as  $\mathbb{E}(A)$  and still using Lemma 12 in [11], we get

$$\begin{aligned} \mathbb{E}(B) &\leq \frac{C}{n - m_n} \sum_{j=k_n+1}^{+\infty} \lambda_j^2 \left( \sum_{j'=1}^{k_n} \frac{\sqrt{\lambda_{j'}}}{\lambda_{j'} - \lambda_j} \Theta v_{j'} \right)^2 \\ &\leq \frac{C}{n - m_n} \sum_{j=k_n+1}^{+\infty} \lambda_j \left( \sum_{j'=1}^{k_n} \frac{\lambda_{j'}}{\lambda_{j'} - \lambda_j} \Theta v_{j'} \right)^2 \\ &\leq \frac{C}{n - m_n} \sum_{j=k_n+1}^{+\infty} \lambda_j \left( \frac{j}{j - k_n} \right)^2 \left( \sum_{j'=1}^{k_n} \Theta v_{j'} \right)^2. \end{aligned}$$

Now, again with the integrability of  $\theta$  and the fact that

$$\sum_{j=k_n+1}^{+\infty} \lambda_j \left( \frac{j}{j - k_n} \right)^2 \leq C k_n b_n,$$

with  $(b_n)_{n \geq 1}$  going to zero as  $n$  goes to infinity (see [11] p.19 in the proof of Proposition 15), we conclude

$$\mathbb{E}(B) \leq \frac{C k_n}{n - m_n} b_n,$$

and this achieves the proof of Lemma 1.

### 5.3 Proof of Lemma 2

Let us denote

$$T_n = \frac{1}{n - m_n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle \delta_i \varepsilon_i.$$

We can write

$$\begin{aligned} T_n^2 &= \frac{1}{(n - m_n)^2} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle^2 \delta_i^2 \varepsilon_i^2 \\ &\quad + \frac{1}{(n - m_n)^2} \sum_{i=1}^n \sum_{\substack{i'=1 \\ i' \neq i}}^n \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle \langle X_{i'}, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle \delta_i \delta_{i'} \varepsilon_i \varepsilon_{i'}. \end{aligned}$$

From the independence between  $\varepsilon$  and  $X$  and the MAR assumption, the expectation of the second term above is zero, hence

$$\begin{aligned}\mathbb{E}(T_n^2) &= \frac{1}{n - m_n} \mathbb{E} \left[ \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle^2 \delta_i^2 \varepsilon_i^2 \right] \\ &= \frac{\sigma_\varepsilon^2}{n - m_n} \mathbb{E} \left[ \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle^2 \delta_i^2 \right],\end{aligned}$$

the index  $i$  corresponding to an observed data in the sample (and consequently  $\delta_i = 1$  for this observation). We finally get

$$\mathbb{E}(T_n^2) = \frac{\sigma_\varepsilon^2}{n - m_n} \mathbb{E} \left[ \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle^2 \right].$$

Following the same lines of the proof of Proposition 17 and Lemma 19 in [11], we obtain

$$\mathbb{E} \left[ \langle X_i, \left( \widehat{\Pi}_{k_n, obs} \widehat{\Gamma}_{n, obs} \right)^{-1} X_\ell \rangle^2 \right] = k_n + o(k_n),$$

which achieves the proof of the Lemma.

#### 5.4 Proof of Lemma 3

The proof of this lemma is quite immediate, noticing that

$$\begin{aligned}\mathbb{E}(\Theta \Pi_{k_n, obs} X_\ell - \Theta X_\ell)^2 &= \mathbb{E}(\langle (\Pi_{k_n, obs} - I) X_\ell, \theta \rangle^2) \\ &= \langle (\Pi_{k_n, obs} - I) \Gamma \theta, \theta \rangle \\ &= \sum_{j=k_n+1}^{+\infty} \left( \Theta \Gamma^{1/2} v_j \right)^2.\end{aligned}$$

#### 5.5 Proof of Theorem 2

From Theorem 1, the last term in the asymptotic development is negligible, so we just have to achieve the usual trade-off between the square bias and the variance. Given that

$$\sum_{j=k_n+1}^{+\infty} \left( \Theta \Gamma^{1/2} v_j \right)^2 = \sum_{j=k_n+1}^{+\infty} L^2 \varphi(j),$$

we approximate this sum with the integral  $\int_x^{+\infty} L^2 \varphi(t) dt$ , which gives the desired result.

## 5.6 Proof of Theorem 3

First, if we follow the same lines of the proof of Lemmas 1 and 3 in Theorem 1 but with all the sample  $X_1, \dots, X_n$ , we get

$$\mathbb{E} \left( \Theta \widehat{\Pi}_{k_n} X_{new} - \Theta \Pi_{k_n} X_{new} \right)^2 = o \left( \frac{\lambda_{k_n} k_n^2}{n} + \frac{k_n}{n} \right), \quad (14)$$

and

$$\mathbb{E} (\Theta \Pi_{k_n} X_{new} - \Theta X_{new})^2 = \sum_{j=k_n+1}^{+\infty} \left( \Theta \Gamma^{1/2} v_j \right)^2. \quad (15)$$

Now, let us denote, for  $i = 1, \dots, n$ ,

$$\varepsilon_{i,imp} = Y_{i,imp} - \langle \theta, X_i \rangle,$$

and

$$\varepsilon_i^* = \delta_i \varepsilon_i + (1 - \delta_i) \varepsilon_{i,imp}.$$

We immediately can write

$$\varepsilon_{i,imp} = \varepsilon_i + Y_{i,imp} - Y_i,$$

and

$$\varepsilon_i^* = \varepsilon_i + (1 - \delta_i)(Y_{i,imp} - Y_i).$$

Then, following the proof of Lemma 2 in Theorem 1, we denote

$$S_n = \frac{1}{n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1} X_{new} \rangle \varepsilon_i^*,$$

whence,

$$\begin{aligned} S_n^2 &= \frac{1}{n^2} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1} X_{new} \rangle^2 (\varepsilon_i^*)^2 \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \sum_{\substack{i'=1 \\ i' \neq i}}^n \langle X_i, \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1} X_{new} \rangle \langle X_{i'}, \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1} X_{new} \rangle \varepsilon_i^* \varepsilon_{i'}^*. \end{aligned}$$

We notice that, for  $i \neq i'$ , we have

$$\mathbb{E} (\varepsilon_i^* \varepsilon_{i'}^*) \leq 4 \mathbb{E} (Y_{i,imp} - Y_i)^2 \leq 8 \left[ \mathbb{E} (Y_{i,imp} - \langle \theta, X_i \rangle)^2 + \sigma_\varepsilon^2 \right].$$

This bound and the lines of the proof of Lemma 2 give

$$\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \langle X_i, \left( \widehat{\Pi}_{k_n} \widehat{\Gamma}_n \right)^{-1} X_{new} \rangle \varepsilon_i^* \right)^2 = O \left( \frac{(n - m_n) k_n}{n^2} + \frac{m_n^2 k_n^2}{n^2} \right). \quad (16)$$

Now, combining relations (14), (15) and (16) and the fact that  $m_n = o(n)$  and  $m_n^2 k_n = O(n)$ , we get the desired result.

## A Appendix

**Table 2** MAR ( $Model_1$ ): Imputed values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)			
Mean	12.520	27.420	44.882
Median	13	27	45
SD	3.307	4.515	5.038
Criterion 1: $[MSE \times 10^3]$	2.3592 (1.8375)	2.7845 (2.0370)	3.2821 (2.0679)
Criterion 2: $[MAE \times 10^2]$	3.7000 (1.4326)	3.9846 (1.4321)	4.3836 (1.3655)
Criterion 3: $[CR3 \times 10^2]$	7.0001 (6.6216)	7.5194 (5.7701)	8.6148 (5.7158)
$n + n_1 = 450$			
Rate of missing data (%)			
Mean	12.433	27.456	45.209
Median	12.333	27.333	45.333
SD	1.877	2.487	3.041
Criterion 1: $[MSE \times 10^3]$	0.8349 (0.5728)	1.0048 (0.6843)	1.3364 (0.9037)
Criterion 2: $[MAE \times 10^2]$	2.2037 (0.7494)	2.4084 (0.8128)	2.7716 (0.9264)
Criterion 3: $[CR3 \times 10^2]$	2.2327 (1.5754)	2.5724 (1.7245)	3.4547 (2.3383)
$n + n_1 = 1800$			
Rate of missing data (%)			
Mean	12.529	27.536	45.213
Median	12.500	27.500	45.250
SD	0.934	1.280	1.355
Criterion 1: $[MSE \times 10^3]$	0.2326 (0.1321)	0.2759 (0.1519)	0.3521 (0.2018)
Criterion 2: $[MAE \times 10^2]$	1.1765 (0.3218)	1.2807 (0.3480)	1.4424 (0.4069)
Criterion 3: $[CR3 \times 10^2]$	0.5933 (0.3492)	0.6962 (0.3891)	0.8822 (0.5036)

**Table 3** MAR ( $Model_2$ ): Imputed values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)			
Mean	12.912	28.026	45.472
Median	13	28	45
SD	3.524	4.493	5.118
Criterion 1: $[MSE \times 10^3]$	2.4786 (2.0871)	2.9537 (2.2814)	3.7448 (2.8036)
Criterion 2: $[MAE \times 10^2]$	3.7528 (1.5474)	4.1059 (1.5388)	4.6319 (1.6688)
Criterion 3: $[CR3 \times 10^2]$	7.5424 (8.0437)	7.7867 (5.7674)	9.7596 (7.1366)
$n + n_1 = 450$			
Rate of missing data (%)			
Mean	12.924	28.018	45.277
Median	13	28	45.33
SD	1.871	2.533	2.844
Criterion 1: $[MSE \times 10^3]$	0.8594 (0.6156)	1.0189 (0.6901)	1.2727 (0.8227)
Criterion 2: $[MAE \times 10^2]$	2.2223 (0.8085)	2.4241 (0.8217)	2.7102 (0.8878)
Criterion 3: $[CR3 \times 10^2]$	2.2861 (1.6605)	2.6008 (1.7465)	3.2415 (2.0856)
$n + n_1 = 1800$			
Rate of missing data (%)			
Mean	13.010	28.081	45.289
Median	13	28.083	45.250
SD	0.970	1.330	1.456
Criterion 1: $[MSE \times 10^2]$	0.1958 (0.1262)	0.2420 (0.1610)	0.2977 (0.1852)
Criterion 2: $[MAE \times 10^2]$	1.0634 (0.34262)	1.1794 (0.3912)	1.3112 (0.4218)
Criterion 3: $[CR3 \times 10^2]$	0.5023 (0.3284)	0.6193 (0.4157)	0.7618 (0.4776)

**Table 4** MCAR ( $Model_1$ ): Imputed values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^3]$	2.3450 (2.0545)	2.7328 (2.0723)	3.1500 (2.0949)
Criterion 2: $[\overline{MAE} \times 10^2]$	3.6740 (1.5563)	3.9815 (1.4156)	4.2895 (1.3849)
Criterion 3: $[\overline{CR3} \times 10^2]$	7.4705 (8.5200)	7.5016 (5.8912)	8.5432 (5.7471)
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^3]$	0.8064 (0.5548)	0.9545 (0.6462)	1.1958 (0.8148)
Criterion 2: $[\overline{MAE} \times 10^2]$	2.1578 (0.7334)	2.3500 (0.7865)	2.6215 (0.8774)
Criterion 3: $[\overline{CR3} \times 10^2]$	2.2137 (1.6687)	2.4524 (1.6200)	3.0869 (2.0933)
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^3]$	0.2233 (0.1260)	0.2608 (0.1469)	0.3182 (0.1769)
Criterion 2: $[\overline{MAE} \times 10^2]$	1.1524 (0.3191)	1.2443 (0.3426)	1.3743 (0.3768)
Criterion 3: $[\overline{CR3} \times 10^2]$	0.5757 (0.3444)	0.6577 (0.3757)	0.8019 (0.4551)

**Table 5** MCAR ( $Model_2$ ): Imputed values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^3]$	2.3450 (2.0545)	2.7328 (2.0723)	3.0946 (2.1618)
Criterion 2: $[\overline{MAE} \times 10^2]$	3.6740 (1.5563)	3.9815 (1.4156)	4.2279 (1.4373)
Criterion 3: $[\overline{CR3} \times 10^2]$	7.4705 (8.5200)	7.5016 (5.8912)	8.3376 (5.8543)
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^3]$	0.7846 (0.5912)	0.8936 (0.6505)	1.0994 (0.7586)
Criterion 2: $[\overline{MAE} \times 10^2]$	2.1094 (0.7779)	2.2476 (0.8129)	2.5089 (0.8770)
Criterion 3: $[\overline{CR3} \times 10^2]$	2.1662 (1.7611)	2.2971 (1.6425)	2.8173 (1.9292)
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
Criterion 1: $[\overline{MSE} \times 10^2]$	0.1904 (0.1246)	0.2312 (0.14623)	0.2847 (0.1697)
Criterion 2: $[\overline{MAE} \times 10^2]$	1.0468 (0.3512)	1.1533 (0.3763)	1.2908 (0.3883)
Criterion 3: $[\overline{CR3} \times 10^2]$	0.4963 (0.3340)	0.5901 (0.3754)	0.7264 (0.4296)



**Table 6** MAR ( $Model_1$ ): Predicted values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

		$n + n_1 = 150$		
Rate of missing data (%)	Mean	12.520	27.420	44.882
	Median	13	27	45
	SD	3.307	4.515	5.038
Criterion 4: $[MSE' \times 10^3]$	2.3383 (1.4987)	2.7173 (1.8390)	3.1939 (2.0391)	
Criterion 5: $[MAE' \times 10^2]$	3.6757 (1.2150)	3.9491 (1.3714)	4.3158 (1.3954)	
Criterion 6: $[CR3' \times 10^2]$	5.9523 (3.7338)	6.9769 (4.9933)	8.2677 (5.6516)	
		$n + n_1 = 450$		
Rate of missing data (%)	Mean	12.433	27.456	45.209
	Median	12.333	27.333	45.333
	SD	1.877	2.487	3.041
Criterion 4: $[MSE' \times 10^3]$	0.8453 (0.5530)	0.9984 (0.6729)	1.3046 (0.8897)	
Criterion 5: $[MAE' \times 10^2]$	2.2109 (0.7064)	2.3926 (0.7893)	2.7376 (0.9182)	
Criterion 6: $[CR3' \times 10^2]$	2.1534 (1.3984)	2.5348 (1.6629)	3.3255 (2.2417)	
		$n + n_1 = 1800$		
Rate of missing data (%)	Mean	12.529	27.536	45.213
	Median	12.500	27.500	45.250
	SD	0.934	1.280	1.355
Criterion 4: $[MSE' \times 10^3]$	0.2295 (0.1282)	0.2746 (0.1512)	0.3474 (0.1982)	
Criterion 5: $[MAE' \times 10^2]$	1.1677 (0.3141)	1.2762 (0.3449)	1.4322 (0.4000)	
Criterion 6: $[CR3' \times 10^2]$	0.5756 (0.3165)	0.6887 (0.3753)	0.8699 (0.4888)	

**Table 7** MAR ( $Model_2$ ): Predicted values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)			
Mean	12.912	28.026	45.472
Median	13	28	45
SD	3.524	4.493	5.118
Criterion 4: $[MSE' \times 10^3]$	2.3556 (1.6157)	2.9148 (2.2111)	3.6204 (2.7093)
Criterion 5: $[MAE' \times 10^2]$	3.6745 (1.2491)	4.0741 (1.4459)	4.5309 (1.6198)
Criterion 6: $[CR3' \times 10^2]$	6.0704 (4.1999)	7.4692 (5.6623)	9.2007 (6.5708)
$n + n_1 = 450$			
Rate of missing data (%)			
Mean	12.924	28.018	45.277
Median	13	28	45.33
SD	1.871	2.533	2.844
Criterion 4: $[MSE' \times 10^3]$	0.8183 (0.5391)	0.9882 (0.6270)	1.2666 (0.8146)
Criterion 5: $[MAE' \times 10^2]$	2.1664 (0.7343)	2.3915 (0.7692)	2.7022 (0.8827)
Criterion 6: $[CR3' \times 10^2]$	2.0977 (1.3686)	2.5322 (1.5836)	3.2364 (2.0620)
$n + n_1 = 1800$			
Rate of missing data (%)			
Mean	13.010	28.081	45.289
Median	13	28.083	45.250
SD	0.970	1.330	1.456
Criterion 4: $[MSE' \times 10^2]$	0.1896 (0.1216)	0.2360 (0.1531)	0.2935 (0.1812)
Criterion 5: $[MAE' \times 10^2]$	1.0461 (0.3391)	1.1647 (0.3851)	1.3029 (0.4177)
Criterion 6: $[CR3' \times 10^2]$	0.4856 (0.3148)	0.6035 (0.3959)	0.7492 (0.4618)

**Table 8** MCAR ( $Model_1$ ): Predicted values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
Criterion 4: $[\overline{MSE'} \times 10^3]$	2.1987 (1.4590)	2.6269 (1.7678)	3.0539 (1.9643)
Criterion 5: $[\overline{MAE'} \times 10^2]$	3.5704 (1.1945)	3.8919 (1.3137)	4.2080 (1.3769)
Criterion 6: $[\overline{CR3'} \times 10^2]$	5.6938 (3.8735)	6.7835 (4.5734)	7.9564 (5.2684)
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
Criterion 4: $[\overline{MSE'} \times 10^3]$	0.8310 (0.5507)	0.9569 (0.6430)	1.1812 (0.8286)
Criterion 5: $[\overline{MAE'} \times 10^2]$	2.1921 (0.7055)	2.3466 (0.7653)	2.6039 (0.8748)
Criterion 6: $[\overline{CR3'} \times 10^2]$	2.11684 (1.3864)	2.4349 (1.5951)	3.0227 (2.1066)
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
Criterion 4: $[\overline{MSE'} \times 10^3]$	0.2229 (0.1237)	0.2620 (0.1496)	0.3184 (0.1787)
Criterion 5: $[\overline{MAE'} \times 10^2]$	1.1506 (0.3094)	1.2451 (0.3433)	1.3732 (0.3768)
Criterion 6: $[\overline{CR3'} \times 10^2]$	0.5589 (0.3056)	0.6579 (0.3721)	0.7986 (0.4458)

**Table 9** MCAR ( $Model_2$ ): Predicted values mean errors and standard deviations for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
Criterion 4: $[MSE' \times 10^3]$	2.1987 (1.4590)	2.6269 (1.768)	3.1007 (2.2980)
Criterion 5: $[MAE' \times 10^2]$	3.5704 (1.2206)	3.8919 (1.2660)	4.2134 (1.4992)
Criterion 6: $[CR3' \times 10^2]$	5.6938 (3.8735)	6.7835 (4.5734)	8.2320 (6.6390)
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
Criterion 4: $[MSE' \times 10^3]$	0.7882 (0.5638)	0.8988 (0.6565)	1.1058 (0.7629)
Criterion 5: $[MAE' \times 10^2]$	2.1154 (0.7590)	2.2558 (0.8192)	2.5095 (0.8682)
Criterion 6: $[CR3' \times 10^2]$	1.9766 (1.3459)	2.2647 (1.5910)	2.7929 (1.8796)
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
Criterion 4: $[MSE' \times 10^2]$	0.1905 (0.1216)	0.2300 (0.1462)	0.2844 (0.1709)
Criterion 5: $[MAE' \times 10^2]$	1.0461 (0.3466)	1.1493 (0.3776)	1.2875 (0.3892)
Criterion 6: $[CR3' \times 10^2]$	0.4843 (.3098)	0.5847 (0.3712)	0.7224 (0.4300)

**Table 10** MAR ( $Model_1$ ): Estimation of  $\theta$  mean square errors, variance and square bias for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)			
Mean	12.520	27.420	44.882
Median	13	27	45
SD	3.307	4.515	5.038
$\overline{MSE''} \times 10^2$	20.33993	22.84329	25.59843
$\overline{Variance} \times 10^2$	16.42143	17.02001	17.58919
$\overline{Bias^2} \times 10^2$	3.918497	5.823277	8.009239
$n + n_1 = 450$			
Rate of missing data (%)			
Mean	12.433	27.456	45.209
Median	12.333	27.333	45.333
SD	1.877	2.487	3.041
$\overline{MSE''} \times 10^2$	8.923099	10.01299	12.37846
$\overline{Variance} \times 10^2$	7.636041	8.680379	10.64885
$\overline{Bias^2} \times 10^2$	1.287058	1.332613	1.729613
$n + n_1 = 1800$			
Rate of missing data (%)			
Mean	12.529	27.536	45.213
Median	12.500	27.500	45.250
SD	0.934	1.280	1.355
$\overline{MSE''} \times 10^2$	3.268755	3.663376	4.294925
$\overline{Variance} \times 10^2$	2.517848	2.870331	3.410527
$\overline{Bias^2} \times 10^2$	0.7509066	0.793045	0.884398

**Table 11** MAR (*Model*<sub>2</sub>): Estimation of  $\theta$  mean square errors, variance and square bias for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)			
Mean	12.912	28.026	45.472
Median	13	28	45
SD	3.524	4.493	5.118
$\overline{MSE''} \times 10^2$	25.77594	30.94147	35.58789
$\overline{Variance} \times 10^2$	17.87099	20.83862	21.5734
$\overline{Biais^2} \times 10^2$	7.904949	10.10285	14.01449
$n + n_1 = 450$			
Rate of missing data (%)			
Mean	12.924	28.018	45.277
Median	13	28	45.33
SD	1.871	2.533	2.844
$\overline{MSE''} \times 10^2$	12.80462	14.15714	16.64587
$\overline{Variance} \times 10^2$	6.696352	8.047992	10.44823
$\overline{Biais^2} \times 10^2$	6.108267	6.109149	6.197638
$n + n_1 = 1800$			
Rate of missing data (%)			
Mean	13.010	28.081	45.289
Median	13	28.083	45.250
SD	0.970	1.330	1.456
$\overline{MSE''} \times 10^2$	7.50709	8.091252	8.477034
$\overline{Variance} \times 10^2$	1.746334	2.096911	2.495418
$\overline{Biais^2} \times 10^2$	5.760756	5.994341	5.981616

**Table 12** MCAR (*Model*<sub>1</sub>): Estimation of  $\theta$  mean square errors, variance and square bias for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	19.94417	22.79952	25.66592
$Variance \times 10^2$	15.47478	17.26436	17.30118
$Biais^2 \times 10^2$	4.469392	5.535164	8.364744
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	8.900453	9.732244	11.44174
$Variance \times 10^2$	7.619475	8.368201	9.842181
$Biais^2 \times 10^2$	1.280978	1.364043	1.599558
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	3.213433	3.567813	4.038663
$Variance \times 10^2$	2.483842	2.781455	3.183366
$Biais^2 \times 10^2$	0.7295905	0.7863581	0.855297

**Table 13** MCAR (*Model*<sub>2</sub>): Estimation of  $\theta$  mean square errors, variance and square bias for samples with different sizes discretized in  $p = 100$  equidistant points based on 500 simulation replications.

$n + n_1 = 150$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	27.05626	29.45794	33.92307
$Variance \times 10^2$	18.66185	20.61093	21.52999
$Biais^2 \times 10^2$	8.394406	8.847012	12.39308
$n + n_1 = 450$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	12.23935	13.62752	16.05732
$Variance \times 10^2$	6.128976	7.515442	9.942408
$Biais^2 \times 10^2$	6.110378	6.112074	6.114911
$n + n_1 = 1800$			
Rate of missing data (%)	10	25	40
$MSE'' \times 10^2$	7.482996	7.923035	8.384754
$Variance \times 10^2$	1.9044	2.106332	2.478208
$Biais^2 \times 10^2$	5.578596	5.816703	5.906546



**Table 14** Real data set: prediction errors over 700 drawn samples.

	$n = 78$ , 8 missing data, 70 observed data		
Test sets	$n/4$	$n/3$	$n/2$
Rate of missing data (%)	13	15	20
$MSE \times 10^2$	24.5650 (8.4750)	25.5172 (8.1444)	29.7827 (15.0889)
$MSA \times 10^2$	37.8834 (6.5104)	38.4424 (6.0372)	41.1194 (8.3055)

**Table 15** MAR ( $Model_1$ ): Computation times and selected dimensions of the CV, GCV and K-fold criteria for samples with different sizes discretized in  $p = 100$  equidistant points.

$n + n_1$	150	450	1800
CV			
Computational times (sec.)	10.5928	74.1095	1158.8180
Best dimension $k_n^*$ (For imputation)	5	5	6
Best dimension $k_n^{**}$ (For prediction)	5	5	6
5-fold CV			
Computational times (sec.)	0.7885	1.3610	4.6047
Best dimension $k_n^*$ (For imputation)	5	5	6
Best dimension $k_n^{**}$ (For prediction)	5	5	6
10-fold CV			
Computational times (sec.)	1.2671	2.6702	9.9181
Best dimension $k_n^*$ (For imputation)	5	5	6
Best dimension $k_n^{**}$ (For prediction)	5	5	6
GCV			
Computational times (sec.)	0.3235	0.4065	1.3558
Best dimension $k_n^*$ (For imputation)	5	5	6
Best dimension $k_n^{**}$ (For prediction)	5	5	6

**Acknowledgements** The authors are grateful to the two anonymous referees for their many valuable comments and constructive suggestions leading to the present version. Thanks to an anonymous associate editor and the journal manager for helpful aspects.

## References

1. Bosq, D., Linear Processes in Function Spaces: Theory and Applications (First edition). NY: Springer, New York (2000).
2. Bugni, F. A., Specification test for missing functional data, *Econometric Theory*, 28, 959–1002 (2012).
3. Cai, T. T. and Hall, P., Prediction in functional linear regression, *The Annals of Statistics*, 34, 2159–2179 (2006).

4. Cardot, H., Ferraty, F. and Sarda, P., Functional linear model, *Statistics and Probability Letters*, 45, 11–22 (1999).
5. Cardot, H., Ferraty, F. and Sarda, P., Spline estimators for the functional linear model, *Statistica Sinica*, 13, 571–591 (2003).
6. Cardot, H. and Johannes, J., Thresholding projection estimators in functional linear models, *Journal of Multivariate Analysis*, 101, 5395–408 (2010).
7. Cheng, P. E., Nonparametric Estimation of Mean Functionals with Data Missing at Random, *Journal of the American Statistical Association*, 89, 81–87 (1994).
8. Chu, C. K. Cheng, P. E., Nonparametric regression estimation with missing data, *Journal of Statistical Planning and Inference*, 48, 85–99 (1995).
9. Chiou, J.-M., Zhang, Y.-C., Chen, W.-H. and Chang, C.-W., A functional data approach to missing value imputation and outlier detection for traffic flow data, *Transportmetrica B: Transport Dynamics*, 2, 106–129 (2014).
10. Crambes, C., Kneip, A. and Sarda, P., Smoothing splines estimators for functional linear regression, *The Annals of Statistics*, 37, 35–72 (2009).
11. Crambes, C. and Mas, A., Asymptotics of prediction in functional linear regression with functional outputs, *Bernoulli*, 19, 2627–2651 (2013).
12. Ferraty, F. and Vieu, P., *Nonparametric functional data analysis: Theory and practice*. NY: Springer-Verlag, New York (2006).
13. Ferraty, F., Sued, M. and Vieu, P., Mean estimation with data missing at random for functional covariables, *Statistics: A Journal of Theoretical and Applied Statistics*, 47, 688–706 (2013).
14. Graham, J. W., *Missing data analysis and design*. NY: Springer, New York (2012).
15. Hall, P. and Horowitz, J. L., Methodology and Convergence Rates for Functional Linear Regression, *The Annals of Statistics*, 35, 70–91 (2007).
16. He, Y., Yucl, R. and Raghunathan, T. E., A functional multiple imputation approach to incomplete longitudinal data, *Statistics in Medicine*, 30, 1137–1156 (2011).
17. Horváth, L. and Kokoszka, P., *Inference for Functional Data with Applications*. NY: Springer-Verlag, New York (2012).
18. Junninen, H., Niska, H., Tuppurainen, K., Ruuskanen, J. and Kolehmainen, M. (2004). Methods for imputation of missing values in air quality datasets. *Atmospheric environment*, 38, 2895–2907.
19. Little, R. J. A. and Rubin, D. B., *Statistical analysis with missing data* (Second edition). NY: John Wiley, New York (2002).
20. Manski, C.F. (1995). *Identification problems in the social sciences*. Harvard University Press.
21. Manski, C.F. (2003). *Partial identification of probability distributions*. Springer-Verlag.
22. Mojirsheibani, M., Nonparametric curve estimation with missing data: A general empirical process approach, *Journal of Statistical Planning and Inference*, 137, 2733–2758 (2007).
23. Preda, C., Saporta, G. and Hadj M. M. H., The NIPALS Algorithm for Functional Data, *Revue Roumaine de Mathématique Pures et Appliquées*, 55, 315–326 (2010).
24. Ramsay, J. O. and Dalzell, C., Some tools for functional data analysis, *Journal Royal Statistical Society B*, 53, 539–572 (1991).
25. Ramsay, J. O. and Silverman, B. W., *Functional Data Analysis* (Second edition). NY: Springer-Verlag, New York (2005).
26. Ramsay, J. O., Hooker, G. and Graves, S., *Functional Data Analysis with R and MATLAB* (First edition). NY: Springer Publishing Company, New York (2009).
27. Shi, J. Q. and Choi, T., *Gaussian Process Regression Analysis for Functional Data*. Chapman and Hall (CRC Press), London (2011).
28. Van Buuren, S., *Flexible Imputation of Missing Data*. NJ: Chapman and Hall (CRC Press), Hoboken (2012).
29. Wang, Q., Linton, O. and Härdle, W., Semiparametric Regression Analysis with Missing Response at Random, *Journal of the American Statistical Association*, 99, 334–345 (2004).
30. Yuan, M. and Cai, T. T., A reproducing kernel Hilbert space approach to functional linear regression, *The Annals of Statistics*, 38, 412–444 (2010).
31. Zhang, J. T., *Analysis of Variance for Functional Data*. NY: Chapman and Hall, New York (2014).