

Prediction in functional regression with discretely observed and noisy covariates

Siegfried Hörmann*

Institute of Statistics, Graz University of Technology
and

Fatima Jammoul

Institute of Statistics, Graz University of Technology

December 14, 2021

Abstract

In practice functional data are sampled on a discrete set of observation points and often susceptible to noise. We consider in this paper the setting where such data are used as explanatory variables in a regression problem. If the primary goal is prediction, we show that the gain by embedding the problem into a scalar-on-function regression is limited. Instead we impose a factor model on the predictors and suggest regressing the response on an appropriate number of factor scores. This approach is shown to be consistent under mild technical assumptions, numerically efficient and gives good practical performance in both simulations as well as real data settings.

Keywords: functional data, factor models, PCA, functional regression, scalar-on-function regression, signal-plus-noise

*Corresponding author. Email: shoermann@tugraz.at

1 Introduction

We consider a sample of functional data $(X_t: 1 \leq t \leq T)$, where each data point X_t corresponds to a curve $(X_t(s): s \in [0, 1])$. These curves can be dependent and stationary or iid. As in practical applications full curves are rarely observed, we furthermore assume that these curves are discretely sampled at the same intraday time points $\mathbf{s} = (s_1, \dots, s_p)$ with $0 \leq s_1 < \dots < s_p \leq 1$. Often these measurements are susceptible to some sort of noise and thus we actually observe

$$Z_t = (X_t(s_1), \dots, X_t(s_p))' + (U_{t1}, \dots, U_{tp})' =: \mathcal{X}_t + U_t. \quad (1)$$

We will assume that the noise vectors (U_t) and signals (\mathcal{X}_t) are independent. We are interested in the situation where p is large and may grow with the sample size. More detailed assumptions on the setting will be given later.

The overall goal of this paper is to investigate a *scalar-on-function* linear model with functional covariates X_t and a scalar response Y_t . Hence we consider the relation

$$Y_t = \alpha + \int_0^1 \beta(s)X_t(s)ds + \varepsilon_t, \quad (2)$$

for some square integrable function β and iid errors (ε_t) . Many contributions have then focused on how to estimate the slope curve $\beta(s)$ and establish consistency, rates of convergence and the like. Arguably the most common estimation approach is based on functional principal components (FPCs). The basic idea is to regress Y_t on the principal component scores of X_t . The resulting coefficients are estimators for the scores of $\beta(s)$, when this function is expanded along the FPCs. For details we refer e.g. to Cardot et al. (1999) or

Hall and Horowitz (2007), where also rates of convergence have been established.

A key difficulty in functional regression is that we are facing an ill-posed problem which requires some regularization techniques. Typical approaches are spectral truncation (e.g. Hörmann and Kidziński (2015)) or Tikhonov regularisation (e.g. Ferraty et al. (2012)). Chakraborty and Panaretos (2019) propose a hybrid version of both approaches. Yuan and Cai (2010) provide a very general method using a reproducing kernel Hilbert space (RKHS). Such regularization techniques can be viewed as different forms of smoothing

and thus it is no surprise that spline-based estimation is another important approach in functional regression. For example, Cardot et al. (2003) used a penalized B-Spline approach reminiscent of ridge regression in order to fit the regression curve.

Most of the papers in FDA literature assume that the explanatory variables X_t are fully observed curves. In practice, however, we hardly ever measure a process over a continuum, but rather face a setting as in (1). Ramsay and Silverman (2005) explore this question from a pragmatic point of view. By smoothing the predictors X as well as the regression function $\beta(s)$ via some set of appropriately chosen basis functions, the problem is cast in a fully functional setup. Li and Hsing (2007) allow the X_t 's to be partially observed as well as potentially noisy. Their estimate for β is again obtained by first smoothing X and then using the result in what is essentially a penalized least squares approach. Kneip et al. (2016) also consider a discrete sampling in a more general regression problem with ‘points of impact’:

$$Y_t = \alpha + \int_0^1 \beta(s)X_t(s)ds + \sum_{r=1}^S \beta_r X_t(\tilde{s}_r) + \varepsilon_t. \quad (3)$$

For improved estimation of this model we refer to Liebl et al. (2020).

In this paper we consider the setting (1). Instead of exploring an estimator for $\beta(s)$ itself, however, we solely focus on the prediction problem. Hence our main target is to find a good predictor

$$\hat{Y}_{T+1} = f_T(Z_1, \dots, Z_T, Z_{T+1}, Y_1, \dots, Y_T).$$

We thus aim to make $\tilde{Y}_{T+1} - \hat{Y}_{T+1}$ small, where $\tilde{Y}_{T+1} = E(Y_{T+1}|X_1, \dots, X_T) = Y_{T+1} - \varepsilon_{T+1}$. Cai and Hall (2006) stress the importance of distinguishing between estimating the regression function and prediction. Is the latter of interest, then regularity of $\hat{\beta}(s)$ may not be the target. They work with fully observed data and derive rates for a fixed non-random regressor x . The results can be extended to noisy data if $n/p = O(1)$. Cardot et al. (2007) and Crambes et al. (2009) consider this problem indirectly. In their setting the predictor functions X_t are observed (potentially with iid noise) on an equidistantly spaced grid. In essence they then study the distance between β and $\hat{\beta}$ in the semi-norm induced by the covariance operator of the X_t . This error in turn is closely related to the error when predicting the conditional mean of Y_{T+1} for any new random function X_{T+1} independent of the sample.

The papers Cai and Hall (2006) and Crambes et al. (2009) are thus probably the closest related to our setup and target and will hence be methods of comparison.

In our theoretical results we will lay particular focus on avoiding restrictive smoothness conditions on the explanatory variables $X_t(s)$ and on the slope function $\beta(s)$. In Section 2 we introduce our method, based on an underlying factor model structure of the noisy covariates. We present our theoretical findings in Section 4. In Section 5 we consider a comprehensive simulation study, showcasing our method in cases of both smooth and highly irregular slope functions. A real data example is given in Section 6. We conclude in Section 7. Proofs and technical lemmas are given in the Appendix.

2 Exploiting the factor model structure

When looking at our regression problem one may wonder if the detour to the functional model (2) is necessary. Since we observe, in fact, a multivariate predictor and not a functional one, it does not seem unreasonable to directly impose a linear model of the form

$$Y_t = \alpha + \sum_{j=1}^p \beta_j X_t(s_j) + \varepsilon_t, \quad (4)$$

and then to explore the problem from a purely multivariate perspective. If the number of observation points p is fixed, and the sampling design is the same for all observations (which is typically the case for machine recorded data), then this approach in principle is doable. In the case of noisy covariates, one must first find a way to eliminate the noise U_t or at least explore how it will impact the inference. However, here we are interested in the setting where p is diverging with the sample size T . In this case the linear regression machinery becomes more delicate, even if the X_t 's were observed without noise. In case p diverges faster than T the problem becomes ill-posed and again requires some regularization approach. Another theoretical and also aesthetic issue is that our model will change with increasing p . Hence the coefficients β_j are actually of the form $\beta_j^{(p)}$. In an asymptotic analysis we may need to specify what the limiting model is, which then naturally brings us back to the functional view in (2). The arguably most important argument against exploiting a model of the form of (4) is the subsequent collinearity issue when estimating the

coefficients β_j . Let $X^{(j)} := (X_1(s_j), \dots, X_T(s_j))'$. If the sampling points \mathbf{s} are dense and if the curves $X_t(s)$ are smooth, then neighboring columns $X^{(j)}$ will resemble each other closely. On the other hand, the smoothness of the curves, which is commonly imposed in the literature, assures that we can very well approximate the linear span of $X = (X_1(\mathbf{s}), \dots, X_T(\mathbf{s}))'$ by a comparably low dimensional subspace. Most common estimation techniques make use of this fact in one way or the other and thus account for the *functional nature of the data*.

In this paper we would like to directly exploit the possibility that X can be sufficiently well approximated by a lower dimensional space. Our approach, however, is not necessarily tied to smoothness. We will explore the *factor space of an approximate factor model* that can be attributed to functional data of the form (1). We begin by explaining how such a factor model and the respective factor space is obtained.

We define the mean function $\mu(s) = EX_t(s)$ and the covariance kernel $\Gamma^X(s, s') = \text{Cov}(X_t(s), X_t(s'))$, respectively. Assuming that Γ^X is continuous, we obtain by the Karhunen-Loève expansion that

$$X_t(s) = \mu(s) + \sum_{\ell \geq 1} x_{t\ell} \varphi_\ell(s), \quad (5)$$

where $\varphi_\ell(s)$ are the eigenfunctions of the covariance operator Γ^X and $x_{t\ell} = \int_0^1 (X_t(s) - \mu(s)) \varphi_\ell(s) ds =: \langle X_t - \mu, \varphi_\ell \rangle$. By Mercer's theorem it follows that the eigenfunctions $\varphi_\ell(s)$ are continuous and that convergence in (5) is uniform, in the sense

$$\sup_{s \in [0,1]} E \left| X_t(s) - \mu(s) - \sum_{\ell=1}^L x_{t\ell} \varphi_\ell(s) \right|^2 \rightarrow 0, \quad L \rightarrow \infty. \quad (6)$$

See e.g. Bosq (2000) for details. The scores ($x_{t\ell}: \ell \geq 1$) are uncorrelated and $\text{Var}(x_{t\ell}) = \lambda_\ell$, where λ_ℓ are the eigenvalues of Γ^X (in decreasing order). Choose some integer $L \geq 1$ and define the matrix

$$B(\mathbf{s}) := (\sqrt{\lambda_1} \varphi_1(\mathbf{s}), \dots, \sqrt{\lambda_L} \varphi_L(\mathbf{s})).$$

Moreover, define $f_t = f_{t,L} = (x_{t1}/\sqrt{\lambda_1}, \dots, x_{tL}/\sqrt{\lambda_L})'$. Then we may write

$$X_t(\mathbf{s}) = \mu(\mathbf{s}) + B(\mathbf{s})f_t + R_t(\mathbf{s}), \quad (7)$$

with $R_t(\mathbf{s}) = R_{t,L}(\mathbf{s}) = X_t(\mathbf{s}) - \mu(\mathbf{s}) - B(\mathbf{s})f_t$ and $\max_j E|R_{t,L}(s_j)| \xrightarrow{P} 0$ as $L \rightarrow \infty$. Hence, $\mu(\mathbf{s}) + B(\mathbf{s})f_t$ provides an explicit form of an L -dimensional proxy of $X_t(\mathbf{s})$.

We are going to incorporate these simple observations in the following manner into our theory.

Assumption 1. *For $L = L(T)$ large enough, we assume that $R_{t,L}(\mathbf{s}) = 0$. The dimension parameter L is allowed to diverge with $T \rightarrow \infty$.*

From a practical point of view, Assumption 1 is not a restriction, since in many real examples $R_{t,L}(s)$ converges to zero rapidly, and hence the error is practically negligible if L is chosen large enough. In its spirit it is related to Assumption (A3) in Crambes et al. (2009), who impose existence of an L dimensional subspace of functions on which $X_t(s)$ can be uniformly sufficiently well approximated on $[0, 1]$. Here we require to have this approximation only on $s \in \{s_1, \dots, s_p\}$, at the price of requiring that the error becomes 0 when L is large enough. Our assumption is a theoretical trade-off, which in turn allows to substitute smoothness assumptions and many additional complex technical constraints which are needed in related papers for deriving theoretical results.

Under (1) and Assumption 1 we then have

$$Z_t = \mu(\mathbf{s}) + B(\mathbf{s})f_t + U_t. \quad (8)$$

It holds that $\text{Var}(f_t) = I_L$ and by assumption $\text{Cov}(f_t, U_t) = 0$, where I_L denotes the identity matrix in $\mathbb{R}^{L \times L}$. Imposing that $\text{Var}(U_t)$ is a diagonal matrix, we see that Z_t follows an L -factor model.

3 Regressing on the factor scores

Let us now assume that we have Model (1) and that Assumption 1 holds. Irrespective of whether we impose the functional linear model (2) or the multivariate linear model (4), using the derived representation for X_t we obtain

$$Y_t = a + b'f_t + \varepsilon_t, \quad (9)$$

where $b = (b_1, \dots, b_L)'$ has components

$$b_\ell = \int_0^1 \beta(s) \sqrt{\lambda_\ell} \varphi_\ell(s) ds \quad \text{or} \quad b_\ell = \sum_{j=1}^p \beta_j \sqrt{\lambda_\ell} \varphi_\ell(s_j),$$

and

$$a = \int_0^1 \mu(s)\beta(s)ds \quad \text{or} \quad a = \sum_{j=1}^p \beta_j \mu(s_j),$$

depending on whether we work under (2) or under (4), respectively. Hence in both cases we obtain an ordinary linear model with explanatory variables f_t , which are however not observable and need to be estimated. For fully observed data, we may use $\hat{f}_t = \int_0^1 X_t(s)\hat{\varphi}_\ell(s)ds/\sqrt{\hat{\lambda}_\ell}$, where $\hat{\lambda}_\ell$ and $\hat{\varphi}_\ell(s)$ denote the empirical eigenvalues and eigenfunctions related to the sample X_1, \dots, X_T . This hence leads to the classical FPC based estimation schemes. In our more realistic setting, however, we don't observe $X_t(s)$ but rather Z_t as in (1). We will thus estimate f_t as the factor scores in a factor model, i.e. we will pursue a purely multivariate scheme instead of a functional one. A specific estimator and its theoretical properties will be discussed in Section 4. Before we go into technical details we summarize our general estimation scheme.

Core algorithm:

1. Estimate $\mu(\mathbf{s})$ by $\hat{\mu}(\mathbf{s}) = \frac{1}{T}(Z_1 + \dots + Z_{T+1})$.
2. Center the data by $\hat{\mu}(\mathbf{s})$.
3. Choose an appropriate order \hat{L} .
4. Compute estimated factor scores \hat{f}_t .
5. Determine \hat{a} and \hat{b} via ordinary least squares.
6. Set $\hat{Y}_{T+1} = \hat{a} + \hat{b}'\hat{f}_{T+1}$.

Some remarks are due.

Remark 1. The factor scores are not unique and can be rotated by an orthogonal matrix $G \in \mathbb{R}^{L \times L}$ leading to an equivalent model. For the purpose of prediction the orientation of \hat{f}_t is irrelevant, because the estimator for the slope will be rotated accordingly. But it has to be noted that the estimators \hat{b} are not directly comparable when the sample size changes. As a consequence of this, the design matrix $\hat{F} = [\hat{f}_1, \dots, \hat{f}_{T+1}]'$ needs to be recalculated whenever the sample size changes, because estimates for the factor f_{T+1} may not follow the same rotation as initial estimates for f_1, \dots, f_T . This implies that

when estimating the factor scores and hence generating the design matrix we need to include the new predictor Z_{T+1} .

Remark 2. An important and non-trivial step in this approach is the estimation of the factor scores \hat{f}_t . In the theoretical framework described in the next section, the estimation is based on a PCA factor model approach. However, the algorithm above can be applied to other factor model methods as well, including a maximum likelihood approach (see e.g. Bai and Li (2012)) or a mixture of PCA and maximum likelihood (see e.g. Bai and Liao (2016)).

Remark 3. A delicate tuning choice in our algorithm concerns the parameter L , which determines the number of factors. A natural approach is to use a cross-validation procedure. We refer e.g. to the work of Owen and Wang (2016) who developed a Bi-Cross-Validation. Onatski (2010) suggests an empirical eigenvalue distribution approach for this problem. In the setting of this paper, we are not specifically interested in the number of factors L that will recover the signal X_t best, but rather the number of factors L that will provide the best predictions. We propose to address this problem by a cross-validation and refer to Section 5 for more details on this.

4 Theoretical bounds for the prediction error

4.1 Assumptions

For a simplified presentation we shall assume from now on that all random variables have zero mean. Fix T for the moment and then define $Z = (Z_1, \dots, Z_{T+1})$ and let $\hat{E} = (\hat{e}_1, \dots, \hat{e}_L)$ be the eigenvectors of $\frac{1}{T+1}Z'Z$ associated with the largest L eigenvalues. We define $\hat{F} = \sqrt{T+1}\hat{E}$, which denotes the proposed estimator for $F = (f_1, \dots, f_{T+1})'$. We note that $\frac{1}{T+1}\hat{F}'\hat{F} = I_L$. We let $Y = (Y_1, \dots, Y_{T+1})'$ and $Y_{(-)} = (Y_1, \dots, Y_T, 0)'$ and define the predictor

$$\hat{Y}_{T+1} := \frac{1}{T+1} \hat{f}'_{T+1} \hat{F}' Y_{(-)}.$$

In other words, the estimate \hat{Y}_{T+1} results from the linear model in which $Y_{(-)}$ has been regressed on \hat{F} . Our primary theoretical goal in this paper is to bound $\hat{Y}_{T+1} - \tilde{Y}_{T+1}$, where $\tilde{Y}_{T+1} = f'_{T+1}b$ is the optimal but infeasible predictor. In the following we list the assumptions we are going to use in our proofs. Assumptions 2–4 stem from Hörmann and Jammoul (2021) and are required to establish the consistent estimation of the factor scores f_t .

Assumption 2. *The noise process (U_t) is i.i.d. zero mean and independent of the signals (X_t) . The processes $(U_{ti}: 1 \leq i \leq p)$ are Gaussian with absolutely summable auto-covariances:*

$$\sum_{h \in \mathbb{Z}} |\gamma^U(h)| \leq C_U < \infty.$$

Assumption 3. *(a) The process $(X_t: t \geq 1)$ is L^4 - m -approximable and has zero mean. (b) The curves $X_t = (X_t(s): s \in [0, 1])$ define fourth order random processes (i.e. $\sup_{s \in [0, 1]} EX_1^4(s) \leq C_X < \infty$) with a continuous covariance kernel.*

Assumption 4. *For the eigenfunctions φ_ℓ of the covariance operator Γ^X it holds that*

$$\max_{1 \leq k, \ell \leq L} \left| \frac{1}{p} \sum_{i=1}^p \varphi_k(s_i) \varphi_\ell(s_i) \right| = O(1)$$

as $T \rightarrow \infty$.

Assumption 5. *The variables (ε_t) are iid, have zero mean and finite variance σ_ε^2 . They are independent of (U_t) and (χ_t) .*

Gaussian errors could be avoided at the expense of requiring certain moment inequalities for the noise processes $(U_{ti}: 1 \leq i \leq p)$. Since we don't require independent noise components, the Gaussian setting is convenient, as the dependence is fully described by the autocovariance function. For ease of presentation, we have chosen to remain within this simplified framework. Furthermore, Assumption 3 shows that we may consider a much broader class of processes in comparison to existing literature. In particular we require no smoothness assumptions on the underlying predictor, aside from a continuous covariance kernel. The notion of L^4 - m -approximability allows for a very general dependence structure between the functional observations, including functional ARMA or functional GARCH models. Assumption 4 is a merely technical assumption. We note that the corresponding sums are proxies for $\int_0^1 \varphi_k(s) \varphi_\ell(s) ds$, which is either zero (when $k \neq \ell$) or one (when $k = \ell$).

4.2 Consistency rates

Now we are ready to formulate our theoretical results. In the first theorem we assume that the order of the factor model L is fixed. We will then increase the complexity of the problem.

Theorem 1. *Consider the functional regression model (2) with sampling scheme (1). Let Assumptions 1–5 hold, where L is fixed. Then*

$$|\widehat{Y}_{T+1} - \widetilde{Y}_{T+1}| = O_P\left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}}\right) \quad \text{as } T \text{ and } p = p(T) \rightarrow \infty. \quad (10)$$

In this setting we hence obtain a parametric rate of convergence, provided that $T/p = O(1)$. Such a rate was also obtained in Cai and Hall (2006) under a quite different setup, invoking a non-random regressor function x with bounded scores $|\langle x, \varphi_k \rangle| \leq Ck^{-\gamma}$ and technical assumptions on γ . As it is pointed out in Crambes et al. (2009), inference on a fixed x cannot, however, be directly compared to rates of convergence of the prediction error for random regressor functions. These authors in turn obtained a non-parametric rate for the mean square prediction error of the order $O(T^{-(2m+2q+1)/(2m+2q+2)})$ where m provides the number of existing derivatives of β and where q is related to the decay rate of λ_j by assuming that $\sum_{j \geq k} \lambda_j = k^{-2q}$. Formally, their result compares to ours with $q \rightarrow \infty$, which would lead again to the same rate as ours.

With a consistent estimator for L , we can obtain the same rates as in Theorem 1.

Corollary 2. *Consider the same setup as in Theorem 1, but assume that L is replaced by a consistent estimator \widehat{L} , then (10) holds.*

We deliberately do not further concretise estimating the number of factors L . This is generally a delicate problem, but in our context, where the focus is on prediction, L can be easily tuned by cross-validation. We stress here that the out-of-sample prediction error is not necessarily minimized when we choose the correct value of L . Let us also note tuning the dimension is a problem which is inherent in other approaches as well. E.g. the consistency rates for the predictor obtained in Cai and Hall (2006) depend on correctly tuning the truncation parameter in the PCA estimator of β . The optimal truncation, leading to the obtained rates, depends on knowledge of the spectrum of Γ^X and the decay rates of the scores of the predictor x and the slope β , when these functions are expanded along the eigenfunctions of Γ^X . In practice, the truncation parameter is also tuned by cross-validation.

Theorem 3. *Let Assumptions 1–5 hold. Assume that $p = p(T) \rightarrow \infty$ and*

$L = L(T) \rightarrow \infty$. Then the one-step prediction error is bounded by

$$|\widehat{Y}_{T+1} - \widetilde{Y}_{T+1}| = O_P \left(\left(\frac{p}{\widehat{\gamma}_L} \right)^4 L^{11/2} \frac{1}{\lambda_L^3} \left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}} \right) \right).$$

Corollary 4. *Suppose that the assumptions of Theorem 3 hold. Furthermore, assume that for some $\nu > 0$ and some $\rho > 0$ we have $\lambda_j \geq \rho j^{-\nu}$ and that there is some $\alpha > 0$ such that $p/\widehat{\gamma}_L = O_P(L^\alpha)$. Then the one-step prediction error is bounded by*

$$|\widehat{Y}_{T+1} - \widetilde{Y}_{T+1}| = O_P \left(L^{4\alpha+3\nu+11/2} \left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}} \right) \right).$$

Polynomial decay rates for the eigenvalues are commonly assumed in related literature. The condition $p/\widehat{\gamma}_L = O_P(L^\alpha)$ is discussed in Hörmann and Jammoul (2021) and can be established if we assume equidistant sampling points and the additional assumption $\sup_s E|X_t(s+h) - X_t(s)|^2 = O(h)$ as $h \rightarrow 0$. The factor $L^{4\alpha+3\nu+11/2}$ may be viewed as a non-parametric convergence rate due to the increase in the dimension of our model. Thus, if L is growing at slow enough polynomial rate we get the convergence in Corollary 4.

5 Simulation study

In order to demonstrate our approach and draw comparisons to common techniques, we present a comprehensive simulation study. To this end, we adapt a set of real data to serve as predictors for the simulation setting. The dataset `pm10` consists of bi-hourly measurements of particulate matter PM10 in Graz, Austria, from October 1st 2010 to March 31st 2011. This dataset has 48 observations over the course of 182 days. To control the smoothness of the underlying signal of our predictor, we pre-smooth this dataset using 21 cubic B-splines. This resulting functional data object is evaluated at $p = 48,96,192$ intraday points. Then, we pull a bootstrap sample of $T = 100, 200, 500, 1000$ curves. The resulting curves represent the underlying signal X_t , for $t = 1, \dots, T$, which have been observed at the equidistant points $s_j \in [0, 1]$ $j = 1, \dots, p$. In a final step, we add iid normal distributed noise U_{tj} to obtain noisy observations $Z_t(s_j) = X_t(s_j) + U_{tj}$. We consider $U_{tj} \sim N(0, \sigma_U^2)$ with two settings $\sigma_U = 2$ and $\sigma_U = 5$. The resulting

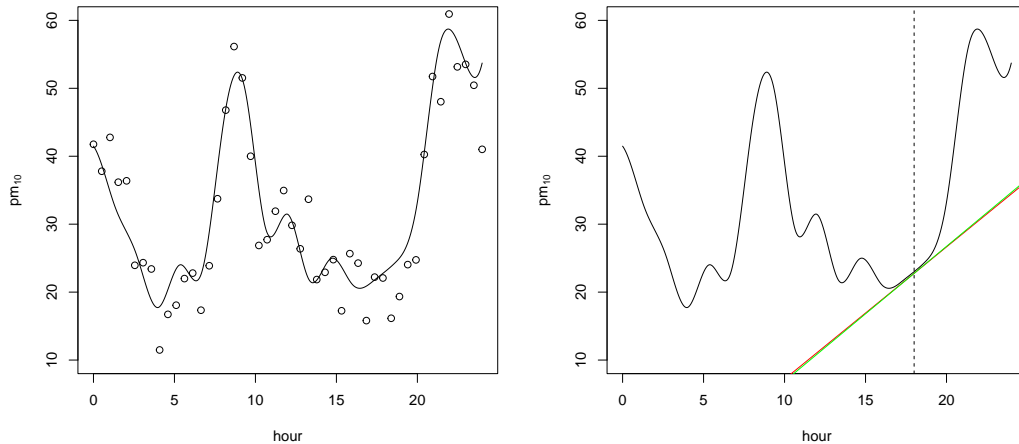


Figure 1: Left: Signal (solid line) and signal-plus-noise with $\sigma_U = 5$ (dots). Right: Slope of $X_t(s)$ the signal at time $x_0 = 18$ (green) and approximation of this slope by $\int_0^{24} \beta_{\text{rough}}(s)X_t(s)ds$ (red).

curves Z_t represent the noisy predictor that we actually observe. One such observation is illustrated in Figure 1.

In the next step, we set up two functional regression models using two very different types of slopes $\beta(s)$. First, we consider the very smooth slope function $\beta_{\text{smooth}}(s) = 10 \sin^3 2\pi s^3$. This particular slope has previously been investigated in Cardot et al. (2007). Additionally, we investigate the effects of a very rough slope

$$\beta_{\text{rough}}(s) := -1/\epsilon^2 I\{s \in [x_0 - \epsilon, x_0)\} + 1/\epsilon^2 I\{s \in [x_0, x_0 + \epsilon]\}.$$

One can see that integration with this slope function is meant to emulate $X'_t(x_0)$ when a small ϵ is chosen. Since keeping the time scale $s \in [0, 1]$ will lead to rather large response values, we have chosen to switch to $s \in [0, 24]$. Then $\int_0^{24} \beta_{\text{rough}}(s)X_t(s)ds \approx X'_t(x_0)$ can be interpreted as the rate of change of the pm10 level at time x_0 in $\mu\text{g}/(\text{m}^3\text{h})$. (See right-hand plot in Figure 1.) In our simulation we choose $\epsilon = 24/100$ and $x_0 = 18$, which corresponds to the tail end of the evening rush hour. In a final step, we obtain the response

Y_t for both slopes by setting

$$Y_t = \int \beta(s)X_t(s)ds + \varepsilon_t$$

where $\varepsilon_t \sim N(0, \sigma_\varepsilon^2)$ with values $\sigma_\varepsilon \in \{2, 5\}$.

Subsequently, we compare three different models: the linear regression model using the recovered factors f_t of Z_t as predictors (**LM**), the smoothing splines approach for noisy predictors as established in Crambes et al. (2009) and Cardot et al. (2007) (**Smooth**) and the functional linear regression with points of impact approach as described in Kneip et al. (2016) and combined with its improved estimation in Liebl et al. (2020) (**PoI**). We note that the method **PoI** is not designed to accommodate noisy observations, and hence the results are not directly compatible with our setup. However, we also experimented with the smoothing splines approach for non-noisy predictors as described in Crambes et al. (2009). These results very closely resembled the **Smooth** approach and hence it is not unreasonable to assume that **PoI** might give competitive results. We also investigated the use of a pre-smoothing step as suggested in Crambes et al. (2009). Here, for more complex noise structures, it was suggested to use a technique to recover the signal from the noisy predictor first and using this estimate to subsequently estimate the functional linear model with Crambes et al. (2009). We found, however, that this extra step didn't give any significant improvement in the predictions compared to the **Smooth** approach, which is why we will not elaborate on these results further.

It is important to note that for the **LM** approach one needs to first estimate the number of factors L required to represent the underlying signal.

Since in this paper we are interested in prediction, we use generalized cross-validation to obtain an estimate for the number of factors. To this end we choose an upper limit L_{\max} and proceed to fit the linear model with all possible choices for $\ell \leq L_{\max}$. We choose the number of factors that minimizes the GCV-score

$$\text{GCV}(\ell) = \frac{T^{-1} \sum_{t=1}^T (Y_t - \widehat{Y}_{t,\ell})^2}{(1 - \ell/T)^2}.$$

In a linear regression this score is a numerically very efficient coefficient which approximates the leave-one-out cross-validation error (see e.g. Hastie et al. (2001)). Note that here, the factors have to only be estimated once using the

maximum number L_{\max} and may then be added one by one into the linear model to then calculate $\text{GCV}(\ell)$, which helps speed up the implementation especially for large datasets. We have set $L_{\max} = 25$.

For the implementation of **PoI** we use the package **FunRegPoI** as provided in the supplementary material to Liebl et al. (2020). This method uses a modified version of the approach described in Crambes et al. (2009) for the estimation of the $\beta(s)$ and β_k . In the smooth setting the number of points of impact is 0. In this case, the estimation corresponds to the method for non-noisy data as described in Crambes et al. (2009). Due to the spiked nature of the slope function β_{rough} we expect that the method might indicate points of impact in this setting.

The implementation of the method **PoI** requires a choice of a maximum number of points of impact S_{\max} , for which we have found the suggested choice $S_{\max} = 8$ to be sufficient. To apply the methods of Crambes et al. (2009) and Liebl et al. (2020), one must initially choose the order of the smoothing splines estimators. In accordance with both references, we have chosen to use cubic smoothing splines. Furthermore, a smoothing parameter must be estimated in the process, which has been achieved via generalized cross validation in analogy to the references.

In Tables 1 and 2 we demonstrate the predictive performances of the different approaches mentioned above. To create these numbers, the dataset was extended by 100 testvalues stemming from the respective model. Then, for the i -th simulation run, we define

$$\text{SSE}_i^{\text{appr}} = \frac{1}{100} \sum_{t=1}^{100} (\tilde{Y}_{T+t} - \hat{Y}_{T+t})^2,$$

where \hat{Y}_t refers to the predicted response. We repeat this simulation for each setting 200 times and report the average of the $\text{SSE}_i^{\text{appr}}$ for $i = 1, \dots, 200$, which we denote by SSE^{appr} . We also report the median for the estimated number of points of impact \hat{S} and the median for the chosen number of factors \hat{L} . We remark that the true L in our chosen setup is 21.

The approaches **Smooth** and **LM** appear to have different strengths when we work with the smooth slope function. The **LM** approach performs best through the setups for large sample sizes and smaller regression errors. In turn, the method **Smooth** works particularly well when we have a small sample size T and larger regression errors. The **PoI** approach is less competitive for this data. It erroneously indicates points of impact. We thus

conjecture that the method is not robust to noisy observations. In the case of β_{rough} (Table 2), the **PoI** can exploit its strength in cases of small p and T but still suffers in other setups. Here **LM** is generally performing best. The most favourable situation for **LM** is when p is small and T is large. It is also of note that in the case of the smooth slope, the number of factors chosen by the GCV tends to be smaller than the true number of factors 21, whereas for the rough slope \hat{L} is close to its actual value.

6 Data illustration

We consider weather-related data recorded in Canadian Weather stations in the neighbouring provinces of Quebec and Ontario. Here, daily mean temperature measurements in the year of 2013 are considered. These yearly curves will serve as our predictors, which we will use to model cumulative log-precipitation in the same year for each station. This type of model has been considered before in Ramsay and Silverman (2005), but we use an extended dataset available at <https://climate.weather.gc.ca/> including more curves.

In an initial data-clean-up phase we eliminate weather stations for which more than 15% of the temperature data is missing. For the remaining missing values, we simply impute them by taking the mean of two adjacent values. If no such mean is available, we simply repeat the last observation. For the associated response, the same removal and imputation methods were used. On occasion, multiple stations with the same name and same geographical position were found in the data. These copies were removed from considerations, as it was difficult to ascertain their nature. After these steps, we have $T = 193$ weather stations with $p = 365$ observations for each predictor curve. A selection of which can be seen in Figure 2. We observe the typical seasonal shape associated with temperature data. The data is then separated randomly into a trainingset of size $T_1 = 143$ and a testset of size $T_2 = 50$. For our analysis the response as well as the predictors have been centered. In a first step, we compute the estimated number of factors using the generalized cross-validation technique. The results of **LM** are subsequently compared to the methods **Smooth** and **PoI**. This split into training- and test set is then repeated 30 times with randomly chosen sets. For the mean sum of squared out-of-sample errors and corresponding standard deviation we obtained 0.0393 (0.0073) for **Smooth**, 0.033 (0.0083)

Dimensions			SSE ^{appr} ($\sigma_U = 2$)					SSE ^{appr} ($\sigma_U = 5$)				
p	T	σ_ε	\hat{L}	\hat{S}	PoI	Smooth	LM	\hat{L}	\hat{S}	PoI	Smooth	LM
48	100	2	15	3	3.63	3.13	3.66	13	2	14.23	13.14	13.91
48	200	2	16	3	2.75	2.5	2.46	14	1	12.07	11.73	11.87
48	500	2	19	3	2.22	2.17	1.97	16	2	10.96	10.94	10.77
48	1000	2	20	4	2.07	2.09	1.84	17	2	10.54	10.6	10.37
96	100	2	15	3	2.92	2.26	2.58	14	1	8.18	7.13	7.85
96	200	2	16	3	1.84	1.62	1.49	15	1	6.44	6.03	6.2
96	500	2	19	3	1.39	1.36	1.1	16	1	5.79	5.55	5.53
96	1000	2	20	4	1.20	1.25	0.94	18	2	5.60	5.47	5.34
192	100	2	15	3	2.28	1.76	1.85	15	1	4.92	3.99	4.68
192	200	2	17	3	1.40	1.2	0.98	16	0	3.81	3.43	3.53
192	500	2	17	4	0.94	0.94	0.62	17	1	3.28	3.1	2.93
192	1000	2	19	4	0.81	0.88	0.53	18	2	3.07	2.97	2.71
48	100	5	13	3	9.64	6.98	9.9	12	1	20.13	16.8	19.25
48	200	5	14	3	5.15	4.02	4.78	13	1	14.70	13.68	14.43
48	500	5	15	3	3.19	2.73	2.95	14	1	11.76	11.53	11.52
48	1000	5	16	3	2.50	2.29	2.3	15	1	11.16	11.03	10.93
96	100	5	13	2	9.48	5.95	8.67	13	1	14.41	10.7	13.86
96	200	5	14	3	4.70	3.16	3.97	14	1	9.10	7.79	8.63
96	500	5	15	2	2.23	1.83	1.98	15	0	6.73	6.25	6.43
96	1000	5	16	3	1.76	1.53	1.39	16	0	6.12	5.88	5.88
192	100	5	13	3	9.05	5.15	7.66	13	1	11.33	7.83	10.74
192	200	5	14	3	4.43	2.76	3.53	14	0	6.53	5.07	5.87
192	500	5	15	3	1.97	1.48	1.55	15	0	4.06	3.66	3.78
192	1000	5	16	3	1.32	1.13	0.96	16	0	3.50	3.24	3.2

Table 1: Simulation Results for the synthetic PM10 data with *smooth* slope function $\beta_{\text{smooth}}(s)$.

Dimensions			SSE ^{appr} ($\sigma_U = 2$)						SSE ^{appr} ($\sigma_U = 5$)					
p	T	σ_ε	\hat{L}	\hat{S}	PoI	Smooth	LM	\hat{L}	\hat{S}	PoI	Smooth	LM		
48	100	2	21	5	7.89	7.81	6.97	19	4	24.85	35.67	28.92		
48	200	2	21	6	6.72	5.66	5.29	21	5	22.04	24.89	22.15		
48	500	2	21	7	5.72	4.51	4.2	21	5	20.75	20.68	18.97		
48	1000	2	21	7	5.36	4.23	3.98	21	6	19.63	19.85	18.27		
96	100	2	21	5	7.87	4.71	4.64	20	4	26.09	19.56	18.68		
96	200	2	21	6	6.43	2.99	2.94	21	5	22.72	13.81	13.61		
96	500	2	21	7	4.96	2.35	2.3	21	6	20.9	11.7	11.49		
96	1000	2	21	8	4.21	2.17	2.12	21	7	18.72	10.93	10.54		
192	100	2	21	6	7.85	3.03	3.09	20	4	25.83	10.19	10.74		
192	200	2	21	6	6.26	1.81	1.83	20	6	21.49	7.56	7.8		
192	500	2	21	7	4.96	1.31	1.29	21	7	18.75	6.29	6.38		
192	1000	2	21	8	4.32	1.16	1.13	21	7	16.72	6.09	6.03		
48	100	5	21	4	12.27	17.29	17.12	18	4	30.52	43.5	37.43		
48	200	5	21	4	8.61	8.69	8.22	20	4	23.87	28.64	25.96		
48	500	5	21	6	6.91	5.75	5.42	21	5	21.7	22.18	20.48		
48	1000	5	21	7	5.88	4.75	4.47	21	6	20.52	20.63	18.88		
96	100	5	20	3	11.5	11.94	12.97	19	4	30.12	28.3	27.14		
96	200	5	21	4	8.52	5.95	6.09	20	5	25.38	16.88	16.89		
96	500	5	21	6	6.18	3.45	3.44	21	5	21.57	12.76	12.5		
96	1000	5	21	7	4.94	2.63	2.59	21	6	20.15	11.65	11.33		
192	100	5	19	3	11.36	10.15	11.06	19	4	30.84	18.31	18.79		
192	200	5	20	4	8.63	4.76	5.04	20	5	24.24	10.47	10.75		
192	500	5	21	6	5.99	2.37	2.4	20	6	19.67	7.29	7.33		
192	1000	5	21	7	5.14	1.65	1.65	21	7	17.85	6.68	6.61		

Table 2: Simulation Results for the synthetic PM10 data with *rough* slope function $\beta_{\text{rough}}(s)$.



Figure 2: Mean Temperature curves in 2013 for 30 out of 193 available Canadian Weather Stations in the Province of Quebec and Ontario

for **PoI** and 0.0307 (0.0084) for **LM**. In comparison, the empirical variance of the log-precipitation is 0.0457 for the total data set. Hence, roughly we can say that our predictions explain approximately 33% (**LM**), 27% (**PoI**) and 14% (**Smooth**) of the variance. We have additionally investigated the difference in the estimated squared out of sample errors in each of the 30 instances using two-sided paired t-tests. The respective p -values are $p = 0.2817$ (comparing **LM** and **PoI**) $p = 8.608 \times 10^{-5}$ (comparing **LM** and **Smooth**). Obviously, the errors between the test sets are not independent and hence most likely the actual differences are less significant. So we cannot conclude to have significantly better predictions from **LM** compared to **PoI**. However, the marked difference between **LM** and **Smooth** presumably cannot be solely attributed to this fact.

To conclude this prediction exercise, we remark that the estimated number of factors \hat{L} was mostly between 26 – 33. The PoI approach estimated between 3 and 5 points of impact, where 2 April (15 times), 25 September (10 times) and 25 May (8 times) were chosen most often.

7 Conclusion

We have been considering scalar-on-function regression with discretely observed and noisy predictors. Although we impose a popular functional model

underlying our data, we can tackle this problem very well from a purely multivariate perspective; this is our key message. The approach we propose is rooted in factor model analysis. It does not require smoothness assumptions on the functional predictors, can accommodate dependent covariates and also allows for dependent sampling errors. We obtain error rates to the theoretically optimal prediction which are comparable to existing results, but work under fewer and simpler assumptions. From a practical side, we show in simulations and real data that our approach gives very convincing results in cases of both smooth and very irregular slope functions. An additional asset is its simple implementation and its fast running time.

A Proofs

We will prove Theorem 3. The proofs of Theorem 1, Corollary 2 and Corollary 4 follow easily from this result.

A.1 Decomposing the prediction error

To expand the prediction error we require further notation. Let us write $B = B(\mathbf{s})$ and introduce

$$H = \frac{1}{T+1} \hat{\Lambda}^{-1} \hat{F}' F B' B,$$

where $\hat{\Lambda} = \text{diag}(\hat{\gamma}_1, \dots, \hat{\gamma}_L)$ with $\hat{\gamma}_1 \geq \dots \geq \hat{\gamma}_L$ being the L largest empirical eigenvalues of $\frac{1}{T+1} Z' Z$. The matrix H takes the role of G in Remark 1. It can be shown to be asymptotically orthogonal and is used to fix an orientation. We write the regression model in the form (9) and then obtain

$$\tilde{Y}_{T+1} = f'_{T+1} H' H b + f'_{T+1} (I_L - H' H) b.$$

Next, let us write (9) in vector notation as

$$Y = F b + \varepsilon,$$

where $\varepsilon = (\varepsilon_1, \dots, \varepsilon_{T+1})'$. From this we obtain the equation

$$b = \frac{1}{T+1} F' Y - \left(\frac{1}{T+1} F' F - I_L \right) b - \frac{1}{T+1} F' \varepsilon.$$

Now we expand the distance between our predictor \widehat{Y}_{T+1} and the best possible predictor \widetilde{Y}_{T+1} in a number of terms which we then shall bound one by one:

$$\begin{aligned}
\widehat{Y}_{T+1} - \widetilde{Y}_{T+1} &= \frac{1}{T+1} \hat{f}'_{T+1} \hat{F}' Y_{(-)} - f'_{T+1} H' H b - f'_{T+1} (I_L - H' H) b \\
&= \frac{1}{T+1} \hat{f}'_{T+1} \hat{F}' Y_{(-)} - \frac{1}{T+1} f'_{T+1} H' H F' Y \\
&\quad + f'_{T+1} H' H \left(\frac{1}{T+1} F' F - I_L \right) b + \frac{1}{T+1} f'_{T+1} H' H F' \varepsilon \\
&\quad + f'_{T+1} (H' H - I_L) b \\
&=: A + B + C,
\end{aligned}$$

where A , B and C correspond to the respective lines of the right hand side of the equation above.

We further expand

$$\begin{aligned}
A &= \frac{1}{T+1} \hat{f}'_{T+1} \hat{F}' Y_{(-)} - \frac{1}{T+1} f'_{T+1} H' H F' Y \\
&= (\hat{f}'_{T+1} - f'_{T+1} H') \frac{1}{T+1} \hat{F}' Y_{(-)} + \frac{1}{T+1} f'_{T+1} H' (\hat{F}' Y_{(-)} - H F' Y) \\
&= (\hat{f}'_{T+1} - f'_{T+1} H') \frac{1}{T+1} \hat{F}' Y_{(-)} \\
&\quad + \frac{1}{T+1} f'_{T+1} H' (\hat{F}' - H F') Y \\
&\quad - \frac{1}{T+1} f'_{T+1} H' \hat{f}_{T+1} Y_{T+1} \\
&:= A_1 + A_2 - A_3,
\end{aligned}$$

where A_1 , A_2 and A_3 correspond to the respective lines of the right hand side of the equation above. Finally, we have

$$B = f'_{T+1} H' H \left(\frac{1}{T+1} F' F - I_L \right) b + \frac{1}{T+1} f'_{T+1} H' H F' \varepsilon := B_1 + B_2.$$

Thus, it remains to investigate A_i and B_i for $i = 1, 2, 3$ and C . These terms are bounded in a series of lemmas in Appendix B and then finally may be combined to obtain the rate given in Theorem 3.

A.2 Technical lemmas

For the reader's convenience, we first summarize some results that are used in the proof of Theorem 3. Lemmas 1–3 constitute slight modifications of results found in Hörmann and Jammoul (2021). In the following, $\|\cdot\|$ denotes the Euclidian norm of a vector or the usual matrix norm, and $\|\cdot\|_F$ is the Frobenius norm.

Lemma 1. *Under Assumptions 1–3 we have*

$$\|\hat{f}_{T+1} - Hf_{T+1}\| = O_P\left(\frac{p}{\hat{\gamma}_L}\sqrt{L}\left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}}\right)\right) \quad (11)$$

and

$$\frac{1}{T+1}\|\hat{F} - H'F'\|_F^2 = O_P\left(L^2\left(\frac{p}{\hat{\gamma}_L}\right)^2\left(\frac{1}{T} + \frac{1}{p}\right)\right). \quad (12)$$

Proof. Statement (12) follows immediately from Lemma 12 in Hörmann and Jammoul (2021). The first statement is a modification of their Lemma 4, in which a uniform bound for $\|\hat{f}_t - Hf_t\|$ over t is derived. Since here we only need a specific value of t , we get the smaller error term. The proof of (11) requires to eliminate some additional factors, which have been used in the proof of Lemma 4 in Hörmann and Jammoul (2021) for bounding the maximum over $t \in \{1, \dots, T+1\}$. This modification actually comes with a similar, but simpler proof and hence the detailed derivations are left to the reader. \square

Lemma 2. *Under Assumptions 1–4 we have*

$$\|H\| = O_P\left(\frac{p}{\hat{\gamma}_L}\frac{L}{\sqrt{\lambda_L}}\right) \quad (13)$$

and

$$\|H'H - I_L\| = O_P\left(\left(\frac{p}{\hat{\gamma}_L}\right)^4\frac{L^5}{\lambda_L^3}\left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}}\right)\right). \quad (14)$$

Proof. This follows immediately from Lemmas 7 and 8 in Hörmann and Jammoul (2021). \square

Lemma 3. *Under Assumption 3 (a) we have*

$$\|(T+1)^{-1}F'F - I_L\|_F = O_P\left(\frac{L}{\lambda_L\sqrt{T}}\right). \quad (15)$$

Proof. This is Lemma 11 in Hörmann and Jammoul (2021). \square

Lemma 4. *We have $\|b\| \leq \sqrt{\lambda_1} \|\beta\|$.*

Proof. We first remark that since β is square integrable, we have by Parseval's identity that $\int_0^1 \beta^2(s) ds = \sum_{k \geq 1} \langle \beta, \varphi_k \rangle^2 < \infty$. Then $\|b\|^2 = \sum_{\ell=1}^L b_\ell^2 \leq \sum_{\ell \geq 1} \lambda_\ell \langle \beta, \varphi_\ell \rangle^2 \leq \lambda_1 \|\beta\|^2$. \square

Lemma 5. *Under Assumptions 2 and 5 we have $\|Y\| = O(\sqrt{T})$.*

Proof. The assumptions imply that (Y_t) defines a stationary and ergodic sequence with 2nd moments. \square

B Proofs

Lemma 6. *Under Assumptions 1–4 we have*

$$\begin{aligned} |A_1| &= O_P \left(\frac{p}{\hat{\gamma}_L} L \left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}} \right) \right); \\ |A_2| &= O_P \left(\left(\frac{p}{\hat{\gamma}_L} \right)^2 \lambda_L^{-1/2} L^{5/2} \left(\frac{1}{T} + \frac{1}{p} \right)^{1/2} \right); \\ |A_3| &= O_P \left(\frac{p}{\hat{\gamma}_L} \lambda_L^{-1/2} L^2 T^{-1/2} \right). \end{aligned}$$

Proof. We have

$$\begin{aligned} |A_1| &= |(f'_{T+1} - f'_{T+1} H') \frac{1}{T+1} \hat{F}' Y_{(-)}| \\ &\leq \sqrt{L} \|f'_{T+1} - f'_{T+1} H'\| \|Y_{(-)}\| / \sqrt{T+1}. \end{aligned}$$

By Lemma 5 we have $\|Y_{(-)}\| / \sqrt{T+1} = O_P(1)$. By (11) the bound for A_1 follows. As for $|A_2|$ we note that

$$\begin{aligned} |A_2| &= \left| \frac{1}{T+1} f'_{T+1} H' (\hat{F}' - H' F') Y \right| \\ &\leq \|f_{T+1}\| \|H\| \frac{1}{\sqrt{T+1}} \left\| \hat{F} - HF \right\|_F \|Y\| / \sqrt{T+1}. \end{aligned}$$

One can easily see that $\|f_{T+1}\| = O_P(\sqrt{L})$. By (12) and (13) the bound for $|A_2|$ follows immediately.

$$\begin{aligned} |A_3| &= \left| \frac{1}{T+1} f'_{T+1} H' \hat{f}_{T+1} Y_{T+1} \right| \\ &\leq \|f_{T+1}\| \|H\| \|\hat{f}_{T+1}/\sqrt{T+1}\| \|Y_{T+1}/\sqrt{T+1}\| \end{aligned}$$

It is easily seen with Markov's inequality that $|Y_{T+1}/\sqrt{T+1}| = O_P(1/\sqrt{T})$. Note furthermore that $\|\hat{f}_{T+1}/\sqrt{T+1}\| \leq \|\hat{F}/\sqrt{T+1}\|_F = \sqrt{L}$. The result follows from (13) and previous considerations. \square

Lemma 7. *Under Assumptions 1–5 we have that*

$$\begin{aligned} |B_1| &= O_P \left(\left(\frac{p}{\hat{\gamma}_L} \right)^2 \lambda_L^{-2} L^{7/2} T^{-1/2} \right); \\ |B_2| &= O_P \left(\left(\frac{p}{\hat{\gamma}_L} \right)^2 \lambda_L^{-1} L^3 T^{-1/2} \right); \end{aligned}$$

Proof. We see that

$$\begin{aligned} |B_1| &= \left| f'_{T+1} H' H \left(\frac{1}{T+1} F' F - I_L \right) b \right| \\ &\leq \|f_{T+1}\| \|H\|^2 \left\| \frac{1}{T+1} F' F - I_L \right\|_F \|b\| \end{aligned}$$

The bound for $|B_1|$ then follows immediately from (13), (15) and Lemma 4.

$$\begin{aligned} |B_2| &= \left| \frac{1}{T+1} f'_{T+1} H' H F' \varepsilon \right| \\ &\leq \|f_{T+1}\| \|H\|^2 \|F' \varepsilon / (T+1)\|. \end{aligned}$$

Using Markov's inequality it is easily seen that $\|F' \varepsilon / (T+1)\| = O_P(\sqrt{L/T})$. The result follows again from (13). \square

Lemma 8. *Under Assumptions 2–4 we have that*

$$|C| = O_P \left(\left(\frac{p}{\hat{\gamma}_L} \right)^4 \lambda_L^{-3} L^{11/2} \left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{p}} \right) \right)$$

Proof. The result follows immediately from the previous considerations and (14). \square

Proof of Theorem 3. Following the Lemmas 6, 7 and 8 we see, given the assumptions of Theorem 3, that the dominant term in the prediction error is $|C|$. The result follows immediately. \square

References

- J. Bai and K. Li. Statistical analysis of factor models of high dimension. *The Annals of Statistics*, 40:436–465, 2012. doi: 10.1214/11-AOS966.
- J. Bai and Y. Liao. Efficient estimation of approximate factor models via regularized maximum likelihood. *Journal of Econometrics*, 191:1–18, 2016.
- D. Bosq. *Linear processes in function spaces: theory and applications*. Lecture Notes in Statistics. Springer, New York, 2000.
- T. Cai and P. Hall. Prediction in functional linear regression. *Ann. Statist.*, 34(5):2159–2179, 10 2006.
- H. Cardot, F. Ferraty, and P. Sarda. Functional linear model. *Statistics & Probability Letters*, 45(1):11–22, 1999.
- H. Cardot, F. Ferraty, and P. Sarda. Spline estimators for the functional linear model. *Statistica Sinica*, 13:571–591, 07 2003.
- H. Cardot, A. Kneip, and P. Sarda. Smoothing spline estimators in functional linear regression with errors-in-variables. *Computational Statistics & Data Analysis*, 51:4832–4848, 06 2007.
- A. Chakraborty and V. Panaretos. Hybrid regularisation and the (in)admissibility of ridge regression in infinite dimensional hilbert spaces. *Bernoulli*, 25:1939–1976, 2019.
- C. Crambes, A. Kneip, and P. Sarda. Smoothing splines estimators for functional linear regression. *Ann. Statist.*, 37(1):35–72, 02 2009.
- F. Ferraty, W. González-Manteiga, A. Martínez-Calvo, and P. Vieu. Presmoothing in functional linear regression. *Statistica Sinica*, 22(1):69–94, 2012.

- P. Hall and J. Horowitz. Methodology and convergence rates for functional linear regression. *Ann. Statist.*, 35(1):70–91, 02 2007.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer, New York, 2001.
- S. Hörmann and F. Jammoul. Consistently recovering the signal from noisy functional data. *Journal of Multivariate Analysis*, 2021. ISSN 0047-259X.
- S. Hörmann and Ł. Kidziński. A note on estimation in hilbertian linear models. *Scandinavian journal of statistics*, 42(1):43–62, 2015.
- A. Kneip, D. Poß, and P. Sarda. Functional linear regression with points of impact. *The Annals of Statistics*, 44(1):1 – 30, 2016.
- Y. Li and T. Hsing. On rates of convergence in functional linear regression. *Journal of Multivariate Analysis*, 98:1782–1804, 10 2007.
- D. Liebl, S. Rameseder, and C. Rust. Improving estimation in functional linear regression with points of impact: Insights into google adwords. *Journal of Computational and Graphical Statistics*, 29(4):814–826, 2020.
- A Onatski. Determining the number of factors from empirical distribution of eigenvalues. *The Review of Economics and Statistics*, 92:1004–1016, 2010.
- A. Owen and J. Wang. Bi-cross-validation for factor analysis. *Statistical Science*, 31:119–139, 2016.
- J.O. Ramsay and B.W. Silverman. *Functional Data Analysis*. Springer, New York, 2005.
- M. Yuan and T. Cai. A reproducing kernel hilbert space approach to functional linear regression. *Ann. Statist.*, 38(6):3412–3444, 12 2010.