

Statistica Sinica Preprint No: SS-2021-0223

Title	Factor-augmented Model for Functional Data
Manuscript ID	SS-2021-0223
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202021.0223
Complete List of Authors	Yuan Gao, Han Lin Shang and Yanrong Yang
Corresponding Authors	Yanrong Yang
E-mails	yanrong.yang@anu.edu.au

FACTOR-AUGMENTED MODEL FOR FUNCTIONAL DATA

Yuan Gao[†], Han Lin Shang^{*} and Yanrong Yang[†]

†Australian National University

** Macquarie University*

Abstract:

We propose modeling raw functional data as a mixture of a smooth function and a high-dimensional factor component. The conventional approach to retrieving a smooth function from raw data is to use a smoothing technique. However, the smoothing model is unable to recover the smooth curve or capture the data variation in some situations, for example, when there is a large measurement error, the smoothing basis functions are incorrectly identified, or the step jumps in the functional mean levels are neglected. We propose a factor-augmented smoothing model to address these challenges, and implement an iterative numerical estimation approach. Including the factor model component in the proposed method solves the aforementioned problems because a few common factors often drive the variation that cannot be captured by the smoothing model. We also establish asymptotic theorems to demonstrate the effects of including factor structures on the smoothing results. Specifically, we show that the smoothing coefficients projected on the complement space of the factor loading matrix are asymptotically normal. Of independent interest, we present an estimator for the population covariance matrix of the raw data, based on the proposed model. Extensive simulation studies show that these factor adjustments are essential to improving the estimation accuracy and avoiding the curse of dimensionality. Lastly, we demonstrate the performance of our model by applying it to Australian

The authors acknowledge editor, associate editor, and two reviewers for their helpful and constructive comments.

temperature data.

Key words and phrases: Basis function misspecification, functional data smoothing, high-dimensional factor model, measurement error, statistical inference on covariance estimation

1. Introduction

Functional data analysis (FDA) has become increasingly popular as data storage technology has improved. Functional data are realizations of smooth random objects using curves, images, and shapes. Ramsay and Silverman (2002, 2005) and Ramsay and Hooker (2017) discuss the methodology and applications of FDA; see also Ferraty and Vieu (2006), Horváth and Kokoszka (2012) Cuevas (2014), Febrero-Bande et al. (2017), Goia and Vieu (2016), Reiss et al. (2017) and Wang et al. (2016). One of the main challenges in an FDA is that we cannot observe functional curves directly, but only discrete points, which are often contaminated by measurement errors. To model a mixture of functional data and high-dimensional measurement errors, we introduce a factor-augmented smoothing model (FASM).

We denote a random sample of n functional data as $\mathcal{X}_i(u)$, for $i = 1, \dots, n$, and $u \in \mathcal{I} \subset \mathbb{R}$, where \mathcal{I} is a compact interval on the real line \mathbb{R} . In practice, the observed data are discrete points and are often contaminated by noise or measurement errors. We use Y_{ij} to represent the j th observation on the i th subject. The observed data can then be expressed as a “signal plus noise” model:

$$Y_{ij} = \mathcal{X}_i(u_j) + \eta_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, p.$$

We use $\mathcal{X}_i(u_j)$ to denote the realization of the j th discrete point on the curve $\mathcal{X}_i(\cdot)$, and η_{ij} is

the noise or measurement error. We assume that a measurement error occurs only where the measurements are taken; thus, the error $\boldsymbol{\eta}_i = (\eta_{i1}, \dots, \eta_{ip})$ is a multivariate term of dimension p . However, in practice, although the signal function component $\boldsymbol{\mathcal{X}}_i = (\mathcal{X}_i(u_1), \dots, \mathcal{X}_i(u_j))$ is of the same dimension p , it differs from $\boldsymbol{\eta}_i$ in nature. Although functions are potentially infinite-dimensional, we may impose smoothing assumptions on them, which usually implies that they possess one or more derivatives. This smoothness feature is used to separate the functions from the measurement errors, that is, a functional smoothing procedure.

When the variance of the noise level is a tiny fraction of the variance of the function, we say the signal-to-noise ratio is high. In this case, classic smoothing tools can be applied to the functional data, including kernel methods (e.g., [Wand and Jones, 1995](#)), local polynomial smoothing (e.g., [Fan and Gijbels, 1996](#)), and spline smoothing (e.g., [Wahba, 1990](#); [Eubank, 1999](#); [Green and Silverman, 1999](#)). Once we have pre-smoothed functions, we can obtain estimates such as the mean and covariance functions. More recent studies on functional smoothing approaches include those of [Cai and Yuan \(2011\)](#), [Yao and Li \(2013\)](#), and [Zhang and Wang \(2016\)](#). In this study, we apply basis smoothing to the functions $\mathcal{X}_i(u)$; that is, we represent $\mathcal{X}_i(u)$ as $\mathcal{X}_i(u) = \sum_{k=1}^K c_{ik} \phi_k(u)$, where $\{\phi_k(u), k = 1, \dots, K\}$ are the basis functions and $\{c_{ik}, i = 1, \dots, n, k = 1, \dots, K\}$ are the smoothing coefficients. The smoothing model then becomes

$$Y_{ij} = \sum_{k=1}^K c_{ik} \phi_k(u_j) + \eta_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, p.$$

When the signal-to-noise level is low, smoothing tools may not adequately remove the measurement error, resulting in an inefficient estimation of the smoothing coefficients. Here we

examine further the measurement error η_{ij} . In an FDA, the number of discrete points p on each subject is often large compared with the sample size n . Hence, the term $\boldsymbol{\eta}_i$ is a high-dimensional component. In this case, the observed data are a mixture of functional data and high-dimensional data. The existence of the large measurement error η_{ij} causes the curse of dimensionality problem, requiring that we apply dimension-reduction models to η_{ij} . Here, factor models are widely used as a dimension-reduction technique for high-dimensional data; among these, factor models are widely used (e.g., [Fan et al., 2008](#); [Lam et al., 2011](#)).

We propose using a factor model for the measurement error term. Without further information on the measurement error, the factor model is appropriate because the estimation of the latent factors does not require any observed variables. The high-dimensional measurement error is assumed to be driven by a small number of unobserved common factors:

$$\eta_{ij} = \mathbf{a}_j^\top \mathbf{f}_i + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, p,$$

where $\mathbf{f}_i \in \mathbb{R}^r$ are the unobserved factors, $\mathbf{a}_j \in \mathbb{R}^r$ are the unobserved factor loadings, r is the number of latent factors, and ϵ_{ij} are idiosyncratic errors with mean zero. Thus, the observed data Y_{ij} can be written as the sum of two components:

$$Y_{ij} = \sum_{k=1}^K c_{ik} \phi_k(u_j) + \mathbf{a}_j^\top \mathbf{f}_i + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, p.$$

This is a basis smoothing model with a factor-augmented form. This proposed model can be modified easily for nonparametric smoothing methods. In [Section S1](#), we illustrate the use of

spline smoothing approaches. In Section S3.3, we apply the nonparametric smoothing model to simulated data.

The following three cases motivate the proposed FASM. In each case, using the proposed model remedies the defects of the traditional smoothing model. Examples of the three cases are provided in Section 2:

1. In traditional smoothing models, the measurement error η_{ij} is assumed to be non-informative and independently and identically distributed (i.i.d.) in both directions. This is an unrealistic assumption when the measurement errors contain information. By applying the factor model, we assume that a small number of unobserved factors capture the covariance in the measurement error. This is usually reasonable in practice, because a few common factors are often the cause of systematic measurement errors.
2. When the smoothing basis functions are incorrectly identified, the smoothing model leads to an erroneous coefficient estimate and large residuals. The proposed model deals with this problem, because the unexplained variation from the misidentification can be modeled using a small number of unobserved common factors.
3. When there are step jumps in the mean levels of the functions, neglecting these mean shifts in a smoothing model results in large residuals at the points where the jumps occur. The changes in the mean levels of the functions come from a universal source, and can be modeled by common factors.

Because the latent factors are unobserved, we propose an iterative approach to simultaneously estimate the smooth function and the factors. We use a principal component analysis (PCA)

to estimate the factor model, and apply a penalized least squares estimation to construct the estimator for the smoothing coefficient c_{ik} . We establish the asymptotic theory of the smoothing coefficient estimator, and prove the consistency of the estimator. We also provide the asymptotic distribution of the projected estimator in the orthogonal complement of the space spanned by the factors \mathbf{f}_i . Here, we see the interplay between the smooth component and the factor model component.

The remainder of this paper is structured as follows. In Section 2, we elaborate on the previously mentioned three motivations, and provide examples. In Section 3, we formally state the model, and provide an iterative estimation approach. We discuss the asymptotic properties of the smoothing coefficients under various assumptions in Section 4. In Section 5, we conduct Monte Carlo simulations on the proposed model under different settings. A real-data example is given in Section 6, and Section 7 concludes the paper. In the Supplementary Material, we provide several important model extensions and data analysis, as well as proofs for the theorems.

2. Motivation

We introduce three examples to motivate the proposed model. In each case, the smoothing model is inadequate for capturing the signal information of the raw data. In the first example, a large measurement error exists, and the residuals after smoothing are large with extreme values. In the second example, the basis functions are selected incorrectly, and part of the functions' variation cannot be captured by the smoothing model. In the third example, there are step jumps in the functional data, and the residuals after smoothing contain gaps. These examples demonstrate

that further modeling of the smoothing residuals is needed.

2.1 Functional data with measurement errors

Figure 1 shows rainbow plots of the average daily temperature and log precipitation at 35 locations in Canada. Because of the nature of the two kinds of data, it is reasonable to assume that temperature and log precipitation are functions over time. However, the two graphs display distinct features. Although there are some perturbations in the temperature plot, it is relatively easy to discern each curve's shape. However, in the precipitation plot, there is a tremendous amount of variability in the raw data, such that it is almost impossible to observe the underlying shape of the curves.

Smooth temperature data can be retrieved without much difficulty using basic smoothing techniques. The residuals are small, with constant variation. On the other hand, the residuals after smoothing exhibit a high level of variation for the precipitation data, and even contain some extreme values. Our model endeavors to explain the large residuals in similar cases to the precipitation data; we show the fitting result in Section S4.

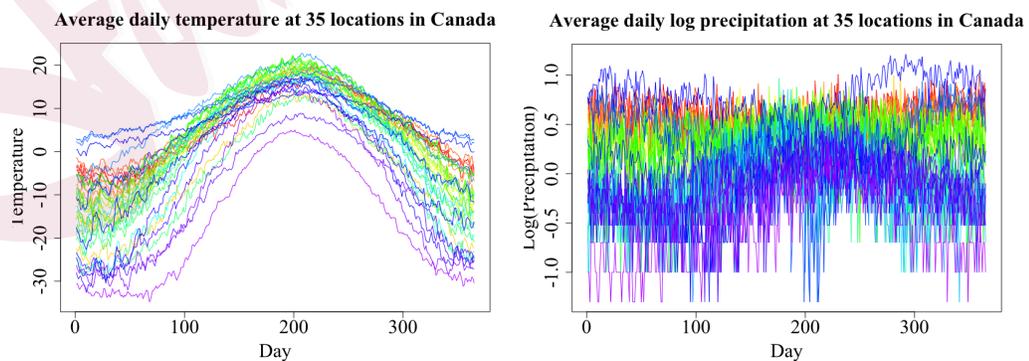


Figure 1: Average daily temperature and log precipitation from 35 Canadian weather stations, averaged over the period 1960 to 1994.

2.2 Misidentification of the basis function

It is important to choose appropriate basis functions for the smoothing method. In this example, we show the inadequacy of the smoothing model when the basis functions are misidentified. We generate functional data using basis functions with changing frequencies. The raw data are shown in Figure 2a. Fourier basis functions are used. In the second half of the data, the frequency of the Fourier basis functions increases, so the data set exhibits more variation toward the right end. Suppose we are unaware of the change in the frequencies of the basis functions, and so used the basis of the first half of the data for the complete curve. The result is shown in Figure 2b. The residuals are large in the second half of the data, but the smoothing model fails to reduce the residuals. Here, a factor model can be used to further model the signal hidden in the large residuals. The data-generating process and further analysis are provided in Section S3.4.

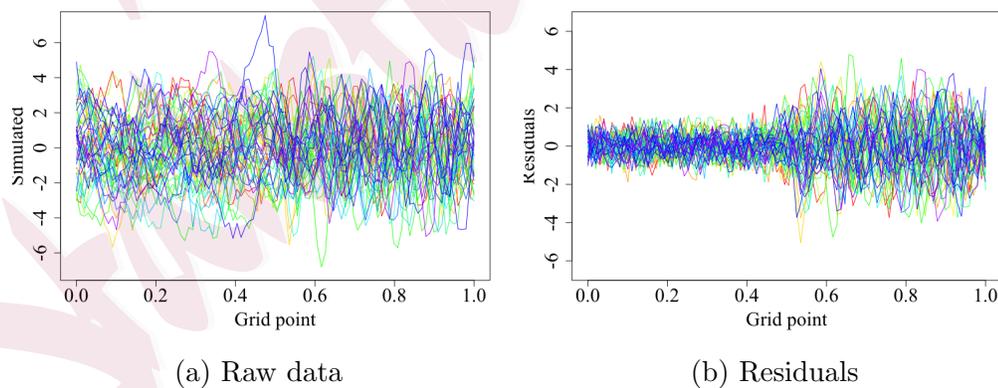


Figure 2: A simulated sample of functional data with changing basis functions.

2.3 Functional data with step jumps in the mean level

Suppose we observe a sample of raw functional data, shown in Figure 3a, where there is a jump at around $u = 0.5$. The jump applies to all sample data, so this sudden shift is at the mean level. We explain how the data are generated in Section S3.5. The residuals after smoothing are presented in Figure 3b. The large residuals around the jump indicate that without measures to deal with the step jumps, smoothing itself is not enough to model these kinds of data. We show in Section S3.5 that applying the proposed model to the same data generates smaller residuals and has less flexibility, which is one of our main goals.

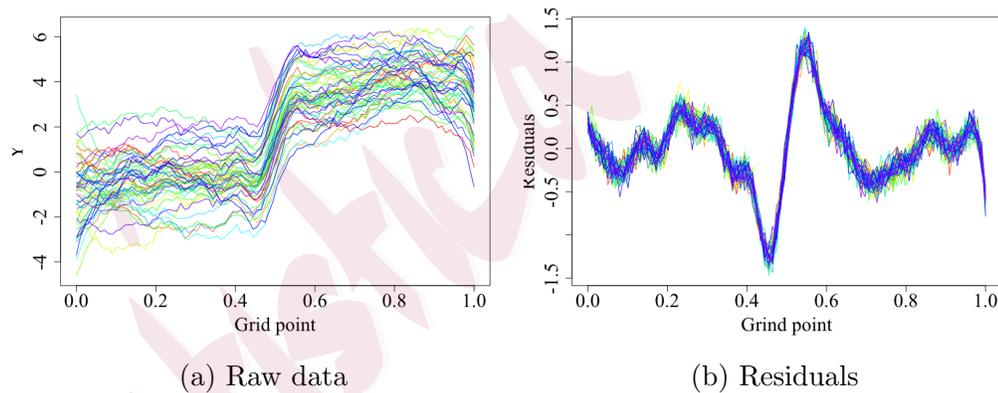


Figure 3: A simulated sample of functional data with step jump.

3. Model Specification and Estimation

We formally state the proposed model in Section 3.1, and present the estimation method in Section 3.2. We first show how the smoothing coefficient c_i and the latent factors \mathbf{f}_i are estimated separately, and then introduce an iterative approach to find these estimates simultaneously.

3.1 Factor-augmented smoothing model

We consider a sample of functional data $\mathcal{X}_i(u)$, that takes values in the space $H := L^2(\mathcal{I})$ of real-valued square integrable functions on \mathcal{I} . The space H is a Hilbert space, equipped with the inner product $\langle x, y \rangle := \int x(u)y(u)du$. The function norm is defined as $\|x\| := \langle x, x \rangle^{1/2}$. The functional nature of $\mathcal{X}_i(u)$ allows us to represent it as a linear expansion of a set of K smooth basis functions:

$$\mathcal{X}_i(u) = \sum_{k=1}^K c_{ik} \phi_k(u), \quad u \in \mathcal{I},$$

where $\{\phi_k(u), k = 1, \dots, K\}$ is a set of common basis functions and c_{ik} is the k th coefficient for the i th curve. Therefore, we can express the full model as

$$Y_{ij} = \sum_{k=1}^K c_{ik} \phi_k(u_j) + \eta_{ij},$$
$$\eta_{ij} = \mathbf{a}_j^\top \mathbf{f}_i + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, p,$$

where $\mathbf{f}_i \in \mathbb{R}^r$ are the unobserved common factors, $\mathbf{a}_j \in \mathbb{R}^r$ are the unobserved factor loadings, and r is the number of factors. We call this model the FASM. For the model to be identifiable, we require the following condition.

Identification Condition 1. We require

(i) $\{\mathcal{X}_i(u_j) : i = 1, \dots, n; j = 1, \dots, p\}$ are independent of $\{\eta_{ij} : i = 1, \dots, n; j = 1, \dots, p\}$;

(ii) $p^{-1} \sum_{j=1}^p \mathbf{a}_j \mathbf{a}_j^\top \xrightarrow{p} \Sigma_{\mathbf{a}} > 0$, for some $r \times r$ matrix $\Sigma_{\mathbf{a}}$ as $p \rightarrow \infty$;

$n^{-1} \sum_{i=1}^n \mathbf{f}_i \mathbf{f}_i^\top \xrightarrow{p} \Sigma_{\mathbf{f}} > 0$, for some $r \times r$ matrix $\Sigma_{\mathbf{f}}$ as $n \rightarrow \infty$.

The first part of the identification condition ensures that the signal function component and the factor model component are independent. The second part ensures the existence of r factors, each of which makes a nontrivial contribution to the variance of η_{ij} , which in turn guarantees the identifiability between the factors and the error term ϵ_{ij} .

We treat the basis functions $\{\phi_k(u) : k = 1, 2, \dots, K\}$ as known, and the number of basis functions K can be fixed or can go to infinity. There are various choices for basis functions in empirical data analysis, and the decision can be subjective. For example, Fourier bases are preferred for periodic data, whereas spline basis systems are most commonly used for nonperiodic data. Other bases include wavelet, polynomial, and ad-hoc basis functions. The number of basis functions K controls the smoothness of the predicted functions. As K increases, the estimator variance increases, but the bias decays. Cross-validation, for instance, can be used to determine K (Wahba and Word, 1975). We estimate our model using a roughness penalty approach explained in the next section. The smoothness of the functional component is switched from being determined by K to being determined by the tuning parameter of the penalty term. In practice, we can use a relatively large K and select the tuning parameter carefully. In Remark 4, we discuss common methods for selecting the tuning parameter.

3.2 Penalized estimation approach

We can write the model for the i th object as

$$\mathbf{Y}_i = \mathbf{\Phi} \mathbf{c}_i + \mathbf{A} \mathbf{f}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, n, \quad (3.1)$$

where

$$\mathbf{Y}_i = \begin{bmatrix} Y_{i1} \\ \vdots \\ Y_{ip} \end{bmatrix}, \mathbf{c}_i = \begin{bmatrix} c_{i1} \\ \vdots \\ c_{iK} \end{bmatrix}, \mathbf{\Phi} = \begin{bmatrix} \phi_1(u_1) & \dots & \phi_K(u_1) \\ \vdots & & \vdots \\ \phi_1(u_p) & \dots & \phi_K(u_p) \end{bmatrix}, \mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_p^\top \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_{i1} \\ \vdots \\ \epsilon_{ip} \end{bmatrix}.$$

Combining all the objects, we have, in matrix form,

$$\mathbf{Y} = \mathbf{\Phi}\mathbf{C} + \mathbf{A}\mathbf{F}^\top + \mathbf{E}, \tag{3.2}$$

where \mathbf{Y} is $p \times n$ and $\mathbf{C} = (\mathbf{c}_1, \dots, \mathbf{c}_n)$ is a $K \times n$ matrix containing all the smoothing coefficients.

The matrix $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_n)^\top$ is $n \times r$ and $\mathbf{E} = (\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_n)$ is $p \times n$. Because $\mathbf{\Phi}$ is assumed to be known, we show how to estimate the parameters \mathbf{C} , \mathbf{A} , and \mathbf{f} .

For the latent factor estimation, there is an identification problem because $\mathbf{A}\mathbf{F}^\top = \mathbf{A}\mathbf{U}\mathbf{U}^{-1}\mathbf{F}^\top$ for any $r \times r$ invertible matrix \mathbf{U} . Thus, we impose the following normalization restriction on the matrices \mathbf{A} and \mathbf{F} :

$$\mathbf{A}^\top \mathbf{A} / p = \mathbf{I}_r, \quad \text{and } \mathbf{F}^\top \mathbf{F} \text{ is a diagonal matrix.} \tag{3.3}$$

We propose implementing a penalized least squares, where the objective function is defined as

$$\text{SSR}(\mathbf{c}_i, \mathbf{A}, \mathbf{f}) = \frac{1}{np} \sum_{i=1}^n [(\mathbf{Y}_i - \mathbf{\Phi}\mathbf{c}_i - \mathbf{A}\mathbf{f}_i)^\top (\mathbf{Y}_i - \mathbf{\Phi}\mathbf{c}_i - \mathbf{A}\mathbf{f}_i) + \alpha \text{PEN}_2(\mathcal{X}_i)],$$

where $\text{PEN}_2(\mathcal{X}_i)$ is a penalty term used for regularization, and α is the tuning parameter con-

trolling the degree of regularization. The same α is used for all the functional observations i . This is a simplified case, where we assume a similar degree of smoothness for all curves. The tuning parameter can be chosen using cross-validation or information criteria. We penalize the “roughness” of the function term. To quantify the notion of “roughness” in a function, we use the square of the second derivative. Define the measure of roughness as

$$\text{PEN}_2(\mathcal{X}_i) = \int_{\mathcal{I}} [D^2 \mathcal{X}_i(s)]^2 ds,$$

where $D^2 \mathcal{X}_i$ denotes taking the second derivative of the function \mathcal{X}_i . Here, a larger tuning parameter α indicates a smoother estimated function. Furthermore, we denote

$$\mathbf{\Phi}(u) = [\phi_1(u), \dots, \phi_K(u)]^\top. \quad (3.4)$$

Then,

$$\mathcal{X}_i(u) = \mathbf{c}_i^\top \mathbf{\Phi}(u).$$

We can re-express the roughness penalty $\text{PEN}_2(\mathcal{X}_i)$ in matrix form, as follows:

$$\begin{aligned}\text{PEN}_2(\mathcal{X}_i) &= \int_{\mathcal{I}} [D^2 \mathcal{X}_i(s)]^2 ds \\ &= \int_{\mathcal{I}} [D^2 \mathbf{c}_i^\top \Phi(s)]^2 ds \\ &= \int_{\mathcal{I}} \mathbf{c}_i^\top D^2 \Phi(s) D^2 \Phi^\top(s) \mathbf{c}_i ds \\ &= \mathbf{c}_i^\top \left[\int_{\mathcal{I}} D^2 \Phi(s) D^2 \Phi^\top(s) ds \right] \mathbf{c}_i \\ &= \mathbf{c}_i^\top \mathbf{R} \mathbf{c}_i, \quad i = 1, \dots, n,\end{aligned}$$

where

$$\mathbf{R} \equiv \int_{\mathcal{I}} D^2 \Phi(s) D^2 \Phi^\top(s) ds. \quad (3.5)$$

The matrix \mathbf{R} is the same for all subjects, and the penalty term $\text{PEN}_2(\mathcal{X}_i)$ differs for each subject only by the coefficient \mathbf{c}_i .

Remark 1. The number of smoothing coefficients \mathbf{c}_i increases as the sample size increases. The penalty term penalizes the “roughness” of the smoothed function and mitigates the effect of the increasing number of parameters, thus controlling the model flexibility.

Therefore, the objective function can be written as

$$\text{SSR}(\mathbf{c}_i, \mathbf{A}, \mathbf{f}) = \frac{1}{np} \sum_{i=1}^n [(\mathbf{Y}_i - \Phi \mathbf{c}_i - \mathbf{A} \mathbf{f}_i)^\top (\mathbf{Y}_i - \Phi \mathbf{c}_i - \mathbf{A} \mathbf{f}_i) + \alpha \mathbf{c}_i^\top \mathbf{R} \mathbf{c}_i],$$

subject to the constraint $\mathbf{A}^\top \mathbf{A}/p = \mathbf{I}_r$.

We aim to estimate the smoothing coefficient \mathbf{c}_i . We left multiply a matrix to each term in (3.1) to project the factor model term onto a zero matrix. Define the projection matrix

$$\mathbf{M}_A \equiv \mathbf{I}_p - \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top = \mathbf{I}_p - \mathbf{A}\mathbf{A}^\top/p. \quad (3.6)$$

Then,

$$\mathbf{M}_A \mathbf{A} \mathbf{f}_i = (\mathbf{I}_p - \mathbf{A}\mathbf{A}^\top/p) \mathbf{A} \mathbf{f}_i = (\mathbf{A} - \mathbf{A}\mathbf{A}^\top \mathbf{A}/p) \mathbf{f}_i = \mathbf{0}.$$

Thjerefore, we estimate \mathbf{c}_i from the projected equation

$$\mathbf{M}_A \mathbf{Y}_i = \mathbf{M}_A \Phi \mathbf{c}_i + \mathbf{M}_A \epsilon_i.$$

The projected objective function becomes

$$\text{SSR}(\mathbf{c}_i, \mathbf{A}) = \frac{1}{np} \sum_{i=1}^n [(\mathbf{M}_A \mathbf{Y}_i - \mathbf{M}_A \Phi \mathbf{c}_i)^\top (\mathbf{M}_A \mathbf{Y}_i - \mathbf{M}_A \Phi \mathbf{c}_i) + \alpha \mathbf{c}_i^\top \mathbf{R} \mathbf{c}_i]. \quad (3.7)$$

By taking the derivative of $\text{SSR}(\mathbf{c}_i, \mathbf{A})$ with respect to each \mathbf{c}_i and setting it to zero, we can solve for the estimator $\hat{\mathbf{c}}_i$:

$$\frac{\partial \text{SSR}(\mathbf{c}_i, \mathbf{A})}{\partial \mathbf{c}_i} = -\frac{1}{np} (\mathbf{M}_A \mathbf{Y}_i - \mathbf{M}_A \Phi \mathbf{c}_i)^\top (\mathbf{M}_A \Phi) + \frac{1}{np} \alpha \mathbf{c}_i^\top \mathbf{R}.$$

Setting the derivative to zero and rearranging the terms, we have

$$(\Phi^\top M_A^\top M_A \Phi + \alpha \mathbf{R}) \mathbf{c}_i = \Phi^\top M_A^\top M_A \mathbf{Y}_i.$$

Using the fact that $M_A^\top M_A = (\mathbf{I}_p - \mathbf{A}\mathbf{A}^\top/p)^\top (\mathbf{I}_p - \mathbf{A}\mathbf{A}^\top/p) = M_A$, we obtain the least squares estimator for \mathbf{c}_i given \mathbf{A} :

$$\hat{\mathbf{c}}_i = (\Phi^\top M_A \Phi + \alpha \mathbf{R})^{-1} \Phi^\top M_A \mathbf{Y}_i.$$

Next, to estimate \mathbf{A} and \mathbf{f}_i , we focus on the factor model

$$\boldsymbol{\eta}_i = \mathbf{A}\mathbf{f}_i + \boldsymbol{\epsilon}_i,$$

expressed in matrix form as

$$\mathbf{Z} = \mathbf{A}\mathbf{F}^\top + \mathbf{E},$$

where $\mathbf{Z} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_n)$. In high-dimensional cases, the unknown factors and loadings are typically estimated using least squares (i.e., the principal component analysis; see, e.g., [Fan et al. \(2008\)](#), [Onatski \(2012\)](#)). The least squares objective function is

$$\text{tr} [(\mathbf{Z} - \mathbf{A}\mathbf{F}^\top)(\mathbf{Z} - \mathbf{A}\mathbf{F}^\top)^\top]. \quad (3.8)$$

Minimizing the objective function with respect to \mathbf{F}^\top , we have $\mathbf{F}^\top = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{Z} = \mathbf{A}^\top \mathbf{Z}/p$,

using (3.3). Substituting in (3.8), we obtain the objective function

$$\begin{aligned} & \text{tr} [(\mathbf{Z} - \mathbf{A}\mathbf{A}^\top \mathbf{Z}/p)(\mathbf{Z} - \mathbf{A}\mathbf{A}^\top \mathbf{Z}/p)^\top] \\ &= \text{tr} (\mathbf{Z}\mathbf{Z}^\top - \mathbf{Z}\mathbf{Z}^\top \mathbf{A}\mathbf{A}^\top/p - \mathbf{Z}\mathbf{Z}^\top \mathbf{A}\mathbf{A}^\top/p + \mathbf{A}\mathbf{A}^\top \mathbf{Z}\mathbf{Z}^\top \mathbf{A}\mathbf{A}^\top/p^2) \\ &= \text{tr}(\mathbf{Z}\mathbf{Z}^\top) - \text{tr}(\mathbf{A}^\top \mathbf{Z}\mathbf{Z}^\top \mathbf{A})/p, \end{aligned}$$

where the last equality uses (3.3) and $\text{tr}(\mathbf{Z}\mathbf{Z}^\top \mathbf{A}\mathbf{A}^\top) = \text{tr}(\mathbf{A}^\top \mathbf{Z}\mathbf{Z}^\top \mathbf{A})$. Thus, minimizing the objective function is equivalent to maximizing $\text{tr}(\mathbf{A}^\top \mathbf{Z}\mathbf{Z}^\top \mathbf{A})/p$. The estimator for \mathbf{A} is obtained by finding the first r eigenvectors corresponding to the r largest eigenvalues of the matrix $\frac{1}{p}\mathbf{Z}\mathbf{Z}^\top$ in descending order, where

$$\frac{1}{p}\mathbf{Z}\mathbf{Z}^\top = \frac{1}{p} \sum_{i=1}^n \boldsymbol{\eta}_i \boldsymbol{\eta}_i^\top = \frac{1}{p} \sum_{i=1}^n (\mathbf{Y}_i - \boldsymbol{\Phi} \mathbf{c}_i)(\mathbf{Y}_i - \boldsymbol{\Phi} \mathbf{c}_i)^\top.$$

Therefore, knowing \mathbf{c}_i , we solve for $\hat{\mathbf{A}}$ using

$$\left[\frac{1}{np} \sum_{i=1}^n (\mathbf{Y}_i - \boldsymbol{\Phi} \mathbf{c}_i)(\mathbf{Y}_i - \boldsymbol{\Phi} \mathbf{c}_i)^\top \right] \hat{\mathbf{A}} = \hat{\mathbf{A}} \mathbf{V}_{np}, \quad (3.9)$$

where \mathbf{V}_{np} is an $r \times r$ diagonal matrix containing the r eigenvalues of the matrix in the square brackets in decreasing order. The additional coefficient $1/n$ is used for scaling.

Remark 2. The number of factors r is selected using some criteria on the eigenvalues. There have been many studies on this topic. For example, Bai and Ng (2002) propose two model selection criteria functions, Onatski (2010) estimate the number of factors using differenced

adjacent eigenvalues, and Ahn and Horenstein (2013) select the number based on the ratio of two adjacent eigenvalues. A modified ratio criterion of Ahn and Horenstein (2013) is proposed in Section S3.1. Note that the simulations and empirical applications presented here show that our model performs well under various common selection criteria in the literature.

We need \mathbf{A} to find $\hat{\mathbf{c}}_i$, and \mathbf{c}_i to find $\hat{\mathbf{A}}$. The final estimator $(\hat{\mathbf{c}}_i, \hat{\mathbf{A}})$ is the solution of the set of equations

$$\begin{cases} \hat{\mathbf{c}}_i = (\Phi^\top M_{\hat{\mathbf{A}}} \Phi + \alpha \mathbf{R})^{-1} \Phi^\top M_{\hat{\mathbf{A}}} \mathbf{Y}_i, & i = 1, \dots, n \\ \left[\frac{1}{np} \sum_{i=1}^n (\mathbf{Y}_i - \Phi \hat{\mathbf{c}}_i) (\mathbf{Y}_i - \Phi \hat{\mathbf{c}}_i)^\top \right] \hat{\mathbf{A}} = \hat{\mathbf{A}} \mathbf{V}_{np}. \end{cases} \quad (3.10)$$

Because there is no closed-form expression of $\hat{\mathbf{A}}$ and $\hat{\mathbf{c}}_i$, we propose using numerical iterations to find the estimates.

Algorithm 1: Iterations for estimating FASM

1. Denote the initial value as $\hat{\mathbf{A}}^{(0)}$. Using (3.10), we obtain $\hat{\mathbf{c}}_i^{(0)} = (\Phi^\top M_{\hat{\mathbf{A}}^{(0)}} \Phi + \alpha \mathbf{R})^{-1} \Phi^\top M_{\hat{\mathbf{A}}^{(0)}} \mathbf{Y}_i$.
 2. With $\hat{\mathbf{c}}_i^{(t)}$, we substitute into the second equation of (3.10) to obtain $\hat{\mathbf{A}}^{(t+1)} = (\hat{\mathbf{a}}_1^{(t+1)}, \dots, \hat{\mathbf{a}}_r^{(t+1)})^\top$, where $\hat{\mathbf{a}}_j^{(t+1)}$ is the eigenvector of the matrix $(np)^{-1} \sum_{i=1}^n (\mathbf{Y}_i - \Phi \hat{\mathbf{c}}_i^{(t)}) (\mathbf{Y}_i - \Phi \hat{\mathbf{c}}_i^{(t)})^\top$ corresponding to its j th largest eigenvalue.
 3. With $\hat{\mathbf{A}}^{(t+1)}$, we obtain $\hat{\mathbf{c}}_i^{(t+1)} = (\Phi^\top M_{\hat{\mathbf{A}}^{(t+1)}} \Phi + \alpha \mathbf{R})^{-1} \Phi^\top M_{\hat{\mathbf{A}}^{(t+1)}} \mathbf{Y}_i$ using (3.10).
 4. We then repeat steps 2 and 3 until $\|\hat{\mathbf{c}}_i^{(t+1)} - \hat{\mathbf{c}}_i^{(t)}\| < \delta$, where δ is a prescribed small positive value.
-

Remark 3. In this paper, we use $\hat{\mathbf{A}}^{(0)} = \mathbf{0}$. This means we start by ignoring the factor model component, so the initial value for the smoothing coefficient $\hat{\mathbf{c}}_i^{(0)} = (\Phi^\top \Phi + \alpha \mathbf{R})^{-1} \Phi^\top \mathbf{Y}_i$, which is simply the ridge estimator. The convergence of Newton's numeric iteration requires the con-

vergence of this estimator, which, in turn, requires the factor model component η_{ij} to have an expectation of zero. The stopping criterion focuses only on $\hat{\mathbf{c}}_i$, because we are interested in estimating η_{ij} as a whole.

Remark 4. Common methods for selecting the shrinkage parameter α include Akaike's information criterion (AIC; Akaike, 1974), the Bayesian information criterion (BIC Schwarz, 1978), and cross-validation. We use the mean generalized cross-validation (mGCV) method (Golub et al., 1979). We define, at step t ,

$$\text{mGCV}^{(t)} = \frac{1}{n} \sum_{i=1}^n \frac{\text{pSSE}_i^{(t)}}{[p - \text{df}^{(t)}(\alpha)]^2}, \quad (3.11)$$

where $\text{SSE}_i^{(t)}$ is the sum of squares residual for the i th object at step t , and $\text{df}^{(t)}(\alpha)$ is the equivalent degrees of freedom measure, which can be calculated as

$$\text{df}^{(t)}(\alpha) = \text{trace} \left[\Phi \left(\Phi^\top \mathbf{M}_{\hat{\mathbf{A}}^{(t)}} \Phi + \alpha \mathbf{R} \right)^{-1} \Phi^\top \mathbf{M}_{\hat{\mathbf{A}}^{(t)}} \right]. \quad (3.12)$$

At each step of the iteration, the tuning parameter α is chosen by minimizing the $\text{mGCV}^{(t)}$.

Remark 5. Algorithm 1 iterates the ridge regression and the PCA. The convergence of this iterative algorithm is studied in Jiang et al. (2021). For instance, Theorem 2 of Jiang et al. (2021) provides sufficient conditions under which the recursive algorithm converges to the true value, or some other values. In particular, this algorithm converges to the true parameter when the regressors are independent of the common factors, or when the factors involved in the regressors

are weaker than the common ones.

After we obtain the estimates $\hat{\mathbf{A}}$ and $\hat{\mathbf{c}}_i$, the estimated coefficient matrix $\hat{\mathbf{C}}$ is constructed as $\hat{\mathbf{C}} = (\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_n)$, and the estimated factor can be obtained using

$$\hat{\mathbf{F}}^\top = \hat{\mathbf{A}}^\top (\mathbf{Y} - \Phi \hat{\mathbf{C}}).$$

Finally, the functional component can be estimated by $\hat{\mathcal{X}}_i(u) = \hat{\mathbf{c}}_i^\top \Phi(u)$, where $\Phi(u)$ is defined in (3.4).

Remark 6. Although we have imposed the constraint in (3.3) and identification condition 1, \mathbf{A} and \mathbf{f}_i are not determined uniquely, because the model (3.1) is unchanged if we replace \mathbf{A} and \mathbf{f}_i with $\mathbf{A}\mathbf{U}$ and $\mathbf{U}^\top \mathbf{f}_i$, for any orthogonal $r \times r$ matrix \mathbf{U} . However, the linear space spanned by the columns of \mathbf{A} is uniquely defined. Although we are not able to estimate \mathbf{A} , we can still estimate a rotation of \mathbf{A} , which spans the same space as \mathbf{A} does. The matrix $\mathbf{M}_{\mathbf{A}}$ defined in (3.6) is a projecting matrix onto the orthogonal complement of the linear space spanned by the columns of \mathbf{A} . In the next section, we show that the estimator $\mathbf{M}_{\hat{\mathbf{A}}}$ for $\mathbf{M}_{\mathbf{A}}$ is consistent.

4. Asymptotic Theory

In this section, we study the asymptotic properties of the coefficient estimator $\hat{\mathbf{c}}_i$ with growing sample size and dimension. We state the assumptions in Section 4.1 and provide the asymptotic results of $\hat{\mathbf{c}}_i$ in Section 4.2.

4.1 Assumptions

We use $(\mathbf{c}_i^0, \mathbf{A}^0)$ to denote the true parameters. In this paper, we use the L^2 norm. The norm of a vector \mathbf{x} is defined as $\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}}$, and the norm of a matrix \mathbf{U} is defined as $\|\mathbf{U}\| = \sqrt{\lambda_{\max}(\mathbf{U}^\top \mathbf{U})}$, where $\lambda_{\max}(\cdot)$ represents the largest eigenvalue of a matrix. We introduce the matrix

$$\mathbf{D}_i(\mathbf{A}) \equiv \frac{1}{p} \Phi^\top \mathbf{M}_A \Phi - \frac{1}{p} \Phi^\top \mathbf{M}_A \Phi \mathbf{f}_i^\top \left(\frac{\mathbf{F}^\top \mathbf{F}}{n} \right)^\top \mathbf{f}_i. \quad (4.1)$$

This matrix plays an important role in the proof of the consistency of $\hat{\mathbf{c}}_i$; see S5.1 Appendix A. The identifying condition for \mathbf{c}_i^0 is that $\mathbf{D}_i(\mathbf{A})$ is positive definite for all i , which is stated in Assumption 3.

First, we state the assumptions.

Assumption 1.

$$\frac{1}{p} \|\Phi^\top \Phi\| = O(1), \quad \|\mathbf{R}\| = O(1), \quad K = o(\min(n, p)), \quad p = o(n^2).$$

In this assumption, the number of basis functions K can either be fixed or tend to infinity.

Assumption 2.

$$\|\mathbf{c}_i^0\| = O_p(1), \quad \text{for all } i = 1, 2, \dots, n.$$

This assumption is introduced to ensure the consistency of the estimated coefficients $\hat{\mathbf{c}}_i^0$. Note that the dimension of \mathbf{c}_i^0 is K , which may go to infinity, but its norm is restricted to be of constant order. Therefore, this assumption controls the flexibility of the model.

Assumption 3. Let $\mathcal{A} = \{\mathbf{A} : \mathbf{A}^\top \mathbf{A}/p = \mathbf{I}, \text{ and } \mathbf{A} \text{ independent of } \Phi\}$. We assume

$$\inf_{\mathbf{A} \in \mathcal{A}} \mathbf{D}_i(\mathbf{A}) > 0.$$

This assumption is the identification condition for \mathbf{c}_i^0 . The usual assumption for the least-squares estimator contains only the first term on the right-hand side of (4.1). The second term arises because of the unobservable matrices \mathbf{F} and \mathbf{A} .

Assumption 4. For some constant $M > 0$,

1. $\mathbb{E}\|\mathbf{a}_j^0\|^4 \leq M$, for $j = 1, \dots, p$, and $\frac{1}{p}\mathbf{A}^\top \mathbf{A} \xrightarrow{p} \Sigma_{\mathbf{a}} > 0$, for some $r \times r$ matrix $\Sigma_{\mathbf{a}}$, as $p \rightarrow \infty$.
2. $\mathbb{E}\|\mathbf{f}_i\|^4 \leq M$ and $\frac{1}{n}\mathbf{F}\mathbf{F}^\top \xrightarrow{p} \Sigma_{\mathbf{f}} > 0$, for some $r \times r$ matrix $\Sigma_{\mathbf{f}}$, as $n \rightarrow \infty$.

Assumption 5. For some constant $M > 0$, the error terms ϵ_{ji} , for $j = 1, \dots, p$ and $i = 1, \dots, n$, are i.i.d. in both directions, with $\mathbb{E}(\epsilon_{ji}) = 0$, $\text{Var}(\epsilon_{ji}) = \sigma^2$, and $\mathbb{E}|\epsilon_{ji}|^8 \leq M$.

Assumption 6. ϵ_{ji} is independent of ϕ_s , \mathbf{f}_t , and \mathbf{a}_s^0 , for all j, i, s, t .

We require that the errors are independent in themselves and also of the functional term $\phi(u)$ and the factor model terms \mathbf{f}_i and \mathbf{a}_j^0 . To highlight the main contribution of our method, we use a simplified setting on the error terms to exclude endogeneity. Nevertheless, with simple, but tedious modifications, Assumption 5 can be relaxed, and our model can be extended to more complicated settings that allow correlations between the error term and the factor model term.

Assumption 7. The tuning parameter α satisfies $\alpha = o(p)$.

This assumption is conventional in ridge regression (see, e.g., Knight and Fu, 2000) and ensures that the estimator's asymptotic bias is zero.

Before stating the next assumption, we introduce some notation. Let ω_j , for $j = 1, \dots, p$, denote the j th column of the $K \times p$ matrix $\Phi^\top \mathbf{M}_{A^0}$, and let ψ_{ik} denote the (i, k) th element of the matrix \mathbf{M}_F , where

$$\mathbf{M}_F \equiv \mathbf{I}_n - \mathbf{F}(\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{F}^\top. \quad (4.2)$$

Then, for any vector $\mathbf{b} = (b_1, \dots, b_n)^\top$, we can write

$$\frac{1}{\sqrt{np}} \Phi^\top \mathbf{M}_{A^0} \mathbf{E} \mathbf{M}_F \mathbf{b} = \frac{1}{\sqrt{np}} \sum_i^n \sum_j^p \omega_j \epsilon_{ji} \sum_k^n \psi_{ik} b_k \equiv \frac{1}{\sqrt{np}} \sum_i^n \sum_j^p \mathbf{x}_{ij}. \quad (4.3)$$

In (4.3), for notational simplicity, we define \mathbf{x}_{ij} as $\omega_j \epsilon_{ji} \sum_k^n \psi_{ik} b_k$. The matrix $\Phi^\top \mathbf{M}_{A^0} \mathbf{E} \mathbf{M}_F$ is of interest because it is the main component that contributes to the asymptotic distribution of the estimators, as shown in the next section.

Let

$$\mathbf{L}_{np} \equiv \frac{\sigma^2}{np} \sum_i^n \sum_j^p \omega_j \omega_j^\top \left(\sum_k^n \psi_{ik} b_k \right)^2. \quad (4.4)$$

We make the following assumption.

Assumption 8. When K is fixed, we assume there exists a $K \times K$ matrix \mathbf{L} such that

$$\mathbf{L} \equiv \lim_{n,p \rightarrow \infty} \mathbf{L}_{np}, \quad (4.5)$$

where \mathbf{L}_{np} is defined in (4.4). Let ν^2 be the smallest eigenvalue of the matrix \mathbf{L} defined in (4.5).

Then assume that $\nu^2 > 0$, and that, for all $\varepsilon > 0$,

$$\lim_{n,p \rightarrow \infty} \frac{1}{np\nu^2} \sum_{i=1}^n \sum_{j=1}^p \mathbb{E} [\|\mathbf{x}_{ij}\|^2 \mathbb{I}(\|\mathbf{x}_{ij}\|^2 \geq \varepsilon np\nu^2)] = 0, \quad (4.6)$$

where $\mathbb{I}(\cdot)$ is the indicator function, defined by $\mathbb{I}(A) = 1$ if A occurs, and $\mathbb{I}(A) = 0$ otherwise.

This assumption is the multivariate Lindeberg condition, which is needed to construct the central limit theorem when K is fixed. This is by no means a strong condition; for instance, when the factor model component is ignored, $\boldsymbol{\omega}_j$ is simply $\boldsymbol{\phi}_j$, and $\mathbf{x}_{ij} = \boldsymbol{\phi}_j b_i \epsilon_{ji}$. Because we assume $\boldsymbol{\phi}_j = O(1)$ in Assumption 1, the Lindeberg condition in (4.6) is met.

Next, we introduce a similar assumption to Assumption 8 for the central limit theorem when K goes to infinity. Let $\tilde{\boldsymbol{\omega}}_j$, for $j = 1, \dots, p$, denote the j th column of the $K \times p$ matrix $\left(\frac{1}{p} \boldsymbol{\Phi}^\top \mathbf{M}_{\mathbf{A}^0} \boldsymbol{\Phi}\right)^{-1} \boldsymbol{\Phi}^\top \mathbf{M}_{\mathbf{A}^0}$. Denote

$$\begin{aligned} & \boldsymbol{\gamma}^\top \left(\frac{1}{p} \boldsymbol{\Phi}^\top \mathbf{M}_{\mathbf{A}^0} \boldsymbol{\Phi}\right)^{-1} \frac{1}{\sqrt{np}} \boldsymbol{\Phi}^\top \mathbf{M}_{\mathbf{A}^0} \mathbf{E} \mathbf{M}_{\mathbf{F}} \mathbf{b} \\ &= \frac{1}{\sqrt{np}} \sum_i^n \sum_j^p \boldsymbol{\gamma}^\top \tilde{\boldsymbol{\omega}}_j \epsilon_{ji} \sum_k^n \psi_{ik} b_k \equiv \frac{1}{\sqrt{np}} \sum_i^n \sum_j^p \tilde{x}_{ij}. \end{aligned} \quad (4.7)$$

Assumption 9. Assume there exists a constant $\tilde{L} > 0$ such that

$$\tilde{L} \equiv \lim_{n,p,K \rightarrow \infty} \tilde{L}_{np}, \quad (4.8)$$

where \tilde{L}_{np} is defined as

$$\tilde{L}_{np} \equiv \frac{\sigma^2}{np} \sum_i^n \sum_j^p (\gamma^\top \tilde{\omega}_j)^2 \left(\sum_k^n \psi_{ik} b_k \right)^2. \quad (4.9)$$

Moreover, for any small value $\varepsilon > 0$, we assume

$$\lim_{n,p,K \rightarrow \infty} \frac{1}{np\nu^2} \sum_{i=1}^n \sum_{j=1}^p \mathbb{E} [|\tilde{x}_{ij}|^2 \mathbb{I}(|\tilde{x}_{ij}|^2 \geq \varepsilon np\nu^2)] = 0, \quad (4.10)$$

where $\mathbb{I}(\cdot)$ is the indicator function, defined by $\mathbb{I}(A) = 1$ if A occurs, and $\mathbb{I}(A) = 0$ otherwise.

4.2 Asymptotic properties

As mentioned previously, the identification problem of the latent factor implies that we use $\hat{\mathbf{A}}$ to estimate a rotation of \mathbf{A}^0 . Based on the objective function (3.7) in Section 3, we use the following center-adjusted objective function:

$$S_{np}(\mathbf{c}_i, \mathbf{A}) = \frac{1}{np} \sum_{i=1}^n [(\mathbf{Y}_i - \Phi \mathbf{c}_i)^\top \mathbf{M}_A (\mathbf{Y}_i - \Phi \mathbf{c}_i) + \alpha \mathbf{c}_i^\top \mathbf{R} \mathbf{c}_i] - \frac{1}{np} \sum_{i=1}^n \epsilon_i^\top \mathbf{M}_{\mathbf{A}^0} \epsilon_i, \quad (4.11)$$

where \mathbf{M}_A is defined in (3.6), satisfying $\mathbf{A}^\top \mathbf{A} / p = \mathbf{I}_r$. The second term on the right-hand side of (4.11) does not contain the unknowns \mathbf{A} and \mathbf{c}_i , so including this term does not affect the optimization result. This term is only used to adjust the center. Thus, the resulting objective

function has an expectation of zero. We estimate \mathbf{c}_i^0 and \mathbf{A}^0 by

$$(\hat{\mathbf{c}}_i, \hat{\mathbf{A}}) = \underset{\mathbf{c}_i, \mathbf{A}}{\operatorname{argmin}} S_{np}(\mathbf{c}_i, \mathbf{A}). \quad (4.12)$$

In the following, we establish the asymptotic properties for the estimated coefficient matrix $\hat{\mathbf{C}}$. In Theorem 1, we prove the consistency of the matrix $\hat{\mathbf{C}}$. In Theorem 2, we show the rate of convergence of $\hat{\mathbf{C}}$. Theorem 3 provides the asymptotic distribution of $\hat{\mathbf{C}}$.

Let $\mathbf{P}_U = \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top$ for a matrix \mathbf{U} .

Theorem 1. *Under Assumptions 1-6, as $n, p \rightarrow \infty$, we have the following statements:*

$$(i) \frac{1}{\sqrt{n}} \|\mathbf{C} - \hat{\mathbf{C}}\| \xrightarrow{p} 0.$$

$$(ii) \|\mathbf{P}_{\hat{\mathbf{A}}} - \mathbf{P}_{\mathbf{A}^0}\| \xrightarrow{p} 0.$$

We start by proving the consistency of the vector $\hat{\mathbf{c}}_i$. This consistency is uniform for all $i = 1, \dots, n$. Therefore, we can combine \mathbf{c}_i , for all $i = 1, \dots, n$, and have the result for the coefficient matrix $\hat{\mathbf{C}}$ in (i). The matrix $\hat{\mathbf{C}}$ is of dimension $K \times n$, where K is fixed and the sample size n goes to infinity, so there is a $(\sqrt{n})^{-1}$ scale in the result of (i). In the second part of the theorem, note that $\mathbf{P}_{\mathbf{A}} = \mathbf{I}_p - \mathbf{M}_{\mathbf{A}}$, where $\mathbf{M}_{\mathbf{A}}$ is the projection matrix onto the orthogonal complement of the linear space spanned by the columns of \mathbf{A} . Thus, $\mathbf{P}_{\hat{\mathbf{A}}}$ and $\mathbf{P}_{\mathbf{A}^0}$ represent the spaces spanned by $\hat{\mathbf{A}}$ and \mathbf{A}^0 , respectively, and we show that they are asymptotically the same in (ii).

Next, we obtain the rate of convergence.

Theorem 2. Under Assumptions 1-6, if $p/n \rightarrow \rho > 0$ as $n, p \rightarrow \infty$, then

$$\left\| \frac{(\mathbf{C}^0 - \widehat{\mathbf{C}})}{\sqrt{n}} \mathbf{M}_{\mathbf{F}} \right\| = O_p \left(\frac{1}{\sqrt{p}} \right),$$

where $\mathbf{M}_{\mathbf{F}}$ is defined in (4.2).

We study the case when the dimension p and the sample size n are comparable. We achieve rate \sqrt{p} convergence, considering $(\sqrt{n})^{-1} \|\mathbf{C}^0 - \widehat{\mathbf{C}}\|$ on the whole. We expect that the rate of convergence of the smoothing models depends on the number of discrete points p observed on each curve.

Remark 7. The asymptotic result in Theorem 2 contains a projection matrix $\mathbf{M}_{\mathbf{F}}$. This matrix projects $\mathbf{C}^0 - \widehat{\mathbf{C}}$ onto the space orthogonal to the factor matrix \mathbf{F} . This theorem shows the relationship between \mathbf{C}^0 and \mathbf{F} . When \mathbf{C}^0 and \mathbf{F} are orthogonal, $(\mathbf{C}^0 - \widehat{\mathbf{C}})\mathbf{M}_{\mathbf{F}} = \mathbf{C}^0 - \widehat{\mathbf{C}}$, and we obtain the rate of convergence of $\mathbf{C}^0 - \widehat{\mathbf{C}}$. When \mathbf{C}^0 and \mathbf{F} are not orthogonal, the inference on \mathbf{C}^0 is affected by the existence of the factor model component.

We further establish the limiting distribution. In S5.1 Appendix A, we show that

$$\left\| \frac{(\mathbf{C}^0 - \widehat{\mathbf{C}})}{\sqrt{p} \sqrt{n}} \mathbf{M}_{\mathbf{F}} \right\| = \left\| \left(\frac{1}{p} \Phi^\top \mathbf{M}_{\mathbf{A}^0} \Phi \right)^{-1} \frac{1}{\sqrt{np}} \Phi^\top \mathbf{M}_{\mathbf{A}^0} \mathbf{E} \mathbf{M}_{\mathbf{F}} \right\| + o_p(1).$$

The limiting distribution is constructed based on the first term on the right-hand side.

Theorem 3. In addition to Assumptions 1-8, we assume that K is fixed and, as $n, p \rightarrow \infty$, there

exists a positive-definite matrix $\tilde{\mathbf{Q}}$ such that

$$\mathbf{Q}(\mathbf{A}^0) \equiv \frac{1}{p} \Phi^\top \mathbf{M}_{\mathbf{A}^0} \Phi \xrightarrow{p} \tilde{\mathbf{Q}}.$$

Then, we have, for any vector $\mathbf{b} \in \mathbb{R}^n$,

$$\sqrt{p} \left(\frac{\mathbf{C}^0 - \hat{\mathbf{C}}}{\sqrt{n}} \right) \mathbf{M}_{\mathbf{F}} \mathbf{b} \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \tilde{\mathbf{Q}}^{-1} \mathbf{L} \tilde{\mathbf{Q}}^{-1} \right),$$

where $\mathbf{M}_{\mathbf{F}}$ is defined in Theorem 2, and \mathbf{L} is defined in (4.5).

Moreover, if the limit of $\Phi \tilde{\mathbf{Q}}^{-1} \mathbf{L} \tilde{\mathbf{Q}}^{-1} \Phi^\top$ exists, we have

$$\sqrt{p} \Phi \left(\frac{\mathbf{C}^0 - \hat{\mathbf{C}}}{\sqrt{n}} \right) \mathbf{M}_{\mathbf{F}} \mathbf{b} \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \lim_{n,p \rightarrow \infty} \Phi \tilde{\mathbf{Q}}^{-1} \mathbf{L} \tilde{\mathbf{Q}}^{-1} \Phi^\top \right),$$

$$\sqrt{\frac{p}{n}} (\mathbf{X} - \hat{\mathbf{X}}) \mathbf{M}_{\mathbf{F}} \mathbf{b} \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \lim_{n,p \rightarrow \infty} \Phi \tilde{\mathbf{Q}}^{-1} \mathbf{L} \tilde{\mathbf{Q}}^{-1} \Phi^\top \right).$$

The vector \mathbf{b} comes from the same vector in Lemma 1. The asymptotic bias is zero, because we assume no serial or cross-sectional correlation in the error terms. This simplified setting can be extended to allow for weak correlations in the errors in both directions. In that case, the asymptotic distribution includes a nonzero bias term.

Remark 8. Theorem 3 shows that the asymptotic distribution of the coefficient matrix $\hat{\mathbf{C}}$ relies on the unobserved factor loading matrix \mathbf{A}^0 . Although we are unable to consistently estimate \mathbf{A}^0 using $\hat{\mathbf{A}}$, what we need is the projected matrix $\mathbf{M}_{\mathbf{A}^0}$, which can be estimated using $\mathbf{M}_{\hat{\mathbf{A}}}$. We

derive the following estimators for \mathbf{Q} and \mathbf{L} based on $\mathbf{M}_{\hat{\mathbf{A}}}$:

$$\begin{aligned}\hat{\mathbf{Q}} &= \frac{1}{p} \Phi^\top \mathbf{M}_{\hat{\mathbf{A}}} \Phi \\ \hat{\mathbf{L}} &= \frac{\sigma^2}{np} \sum_i^n \sum_j^p \hat{\omega}_j^\top \hat{\omega}_j \left(\sum_k^n \hat{\psi}_{ik} b_k \right)^2,\end{aligned}$$

where $\hat{\omega}_j$ is the j th column of the $K \times p$ matrix $\Phi^\top \mathbf{M}_{\hat{\mathbf{A}}}$, and $\hat{\psi}_{ik}$ is the (i, k) th element in the matrix $\mathbf{M}_{\hat{\mathbf{F}}}$.

When K goes to infinity, the random vector $\left(\frac{\mathbf{C}^0 - \hat{\mathbf{C}}}{\sqrt{n}} \right) \mathbf{M}_{\mathbf{F}} \mathbf{b}$ considered in Theorem 3 is a high-dimensional random vector. Then, we establish the asymptotic distribution for some linear combination of this random vector in the following theorem.

Theorem 4. *Suppose Assumptions 1-7 and Assumption 9 are satisfied. As n, p, K tend to infinity, we have, for any vectors $\boldsymbol{\gamma} \in \mathbb{R}^K$ and $\mathbf{b} \in \mathbb{R}^n$,*

$$\sqrt{p} \boldsymbol{\gamma}^\top \left(\frac{\mathbf{C}^0 - \hat{\mathbf{C}}}{\sqrt{n}} \right) \mathbf{M}_{\mathbf{F}} \mathbf{b} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{L}}), \quad (4.13)$$

where $\tilde{\mathbf{L}}$ is defined in (4.8).

5. Simulation Studies

In this section, we use simulated data to demonstrate the performance of the proposed model. We compare the FASM with three other models. The first is the basis smoothing model with a penalty (Bsmooth). We use the same basis functions as those of the FASM and the L_2 -penalty.

Thus, the Bsmooth model differs from the FASM only in that it does not consider the augmented factor component. The second model is the local linear smoothing model (Localin). Here, we use Cross-validation to choose the bandwidth parameter. The third model is the principal component analysis with conditional expectation (PACE), proposed in Yao et al. (2005). This method is mainly used for sparse data, where very few data points are observed on each functional curve.

5.1 Data-generating process

Setting 1

We generate simulated data Y_{ij} , where $i = 1, \dots, n$ and $j = 1, \dots, p$, from the following model:

$$\begin{aligned} Y_{ij} &= \mathcal{X}_i(u_j) + \eta_{ij} + \epsilon_{ji} \\ &= \sum_{k=1}^{13} c_{ik} \phi_k(u_j) + \sum_{k=1}^4 f_{ik} a_{kj} + \epsilon_{ji}, \end{aligned}$$

where $\phi_k(u)$ are chosen as B-spline basis functions of order four and the smoothing coefficients c_{ik} are generated from $\mathcal{N}(0, 1.5^2)$. The factor loadings a_{kj} follow $\mathcal{N}(0, 0.6^2)$, and the factors $(f_{i1}, f_{i2}, f_{i3}, f_{i4})^\top \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is a 4×4 covariance matrix. We set the multivariate mean term $\boldsymbol{\mu} = \mathbf{0}$ and the variance $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_4$. We adjust the value of σ^2 to control the signal-to-noise ratio. When σ^2 is large, the signal-to-noise level is low, and when σ^2 is small, the level is high. The random error terms ϵ_{ji} follow $\mathcal{N}(0, 0.5^2)$.

Setting 2

We generate simulated data Y_{ij} , where $i = 1, \dots, n$ and $j = 1, \dots, p$, from the following model:

$$\begin{aligned} Y_{ij} &= \mathcal{X}_i(u_j) + \eta_{ij} + \epsilon_{ji} \\ &= \sum_{k=1}^5 c_{ik} \phi_k(u_j) + \sum_{k=1}^4 f_{ik} a_{kj} + \epsilon_{ji}, \end{aligned}$$

where $\phi_k(u)$ are chosen as B-spline basis functions of order four and the smoothing coefficients c_{ik} are generated from $\mathcal{N}(0, \gamma_k^2)$, where $\gamma_1 = 3, \gamma_2 = 2.5, \gamma_3 = 2, \gamma_4 = 1.5$, and $\gamma_5 = 1$. The factor loadings a_{kj} follow $\mathcal{N}(0, 0.8^2)$ and the factors $(f_{i1}, f_{i2}, f_{i3}, f_{i4})^\top \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is a 4 by 4 covariance matrix. We set the multivariate mean term $\boldsymbol{\mu} = \mathbf{0}$ and variance $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_4$. We adjust the value of σ^2 to control the signal-to-noise ratio. When σ^2 is large, the signal-to-noise level is low, and when σ^2 is small, the signal-to-noise level is high. The random error terms ϵ_{ji} follow $\mathcal{N}(0, 0.5^2)$.

Setting 3

We generate simulated data as in Setting 1, except that we use $K = 21$ B-spline basis functions. We use $n = 40$ and $p = 50$. Here, we show how the model works when the number of basis functions is large compared with n and p .

5.2 Results

We repeat the simulation setup 500 times and obtain the estimated smooth function $\hat{\mathcal{X}}_i(u) = \hat{\mathbf{c}}_i^\top \Phi(u)$. The mean integrated squared error (MISE) for the function estimation is calculated as

$$\text{MISE} = \frac{1}{n} \sum_{i=1}^n \int [\mathcal{X}_i(u) - \hat{\mathcal{X}}_i(u)]^2 du. \quad (5.1)$$

The results for Settings 1 and 2 are reported in Tables 1 and 2, respectively. With the same sample size n , increasing the number of points p on the curve decreases the estimation error. However, with the same value for p , increasing the sample size does not decrease the estimation error. This is consistent with the convergence rate stated in Section 4, where the estimator converges with the rate related to p .

When σ is large, that is, the signal-to-noise ratio is high, the FASM outperforms the smoothing model. When $\sigma = 0$, the actual data are generated without the factor model component. The performance of the two methods is identical. This implies that the proposed model is robust in that it improves the estimation when there is a measurement error, and is no worse than the ordinary smoothing model when there is no measurement error. This is credited to the selection criterion used in (S3.4), which always selects zero factors when there are no common factors.

In Table 3, we show the MISE result from Setting 3. Again, the proposed model shows its superiority when K is large compared with n and p . This setting corresponds to allowing the number of basis functions K to go to infinity, discussed in Section 4.

Table 1: **Setting 1:** The MISE of the function estimates with different samples sizes and dimensions. The σ^2 -value is used to control the signal-to-noise ratio.

Sample Size	Dimension	Size of η	MISE			
			FASM	Bsmooth	Localin	PACE
$n = 20$	$p = 51$	$\sigma = 0$	0.058	0.058	0.075	0.238
		$\sigma = 0.5$	0.131	0.131	0.149	0.587
		$\sigma = 0.75$	0.204	0.208	0.226	1.026
		$\sigma = 1$	0.297	0.318	0.336	1.676
$n = 20$	$p = 101$	$\sigma = 0$	0.031	0.031	0.043	0.238
		$\sigma = 0.5$	0.076	0.076	0.089	0.595
		$\sigma = 0.75$	0.115	0.123	0.136	1.020
		$\sigma = 1$	0.159	0.187	0.201	1.686
$n = 50$	$p = 51$	$\sigma = 0$	0.058	0.058	0.076	0.195
		$\sigma = 0.5$	0.131	0.135	0.153	0.543
		$\sigma = 0.75$	0.186	0.214	0.233	1.013
		$\sigma = 1$	0.248	0.314	0.331	1.668
$n = 100$	$p = 101$	$\sigma = 0$	0.031	0.031	0.044	0.170
		$\sigma = 0.5$	0.062	0.074	0.089	0.463
		$\sigma = 0.75$	0.090	0.118	0.132	0.908
		$\sigma = 1$	0.133	0.181	0.194	1.570

6. Application to Climatology

Here, we apply the aforementioned FASM, Bsmooth, Localin, and PACE methods to real data sets. This section examines the Friday temperature data at Adelaide airport. We provide analysis of Canadian annual temperature and precipitation data in the Supplementary Material.

We choose Adelaide because it tends to have the hottest temperature among Australia's big cities. Data from other weekdays exhibit similar features, and thus are not shown here. The data are measured every half an hour for the period 1997 to 2007. The sample size n is 508, and the number of discrete data points p from each curve is 48. A plot of the raw data is shown in Figure 4a. The data are quite noisy, with extreme values in some of the curves due to large

Table 2: **Setting 2:** The MISE of the function estimates with different samples sizes and dimensions. The σ^2 -value is used to control the signal-to-noise ratio.

Sample Size	Dimension	Size of η	MISE			
			FASM	Bsmooth	Localin	PACE
$n = 20$	$p = 51$	$\sigma = 0$	0.024	0.024	0.070	0.116
		$\sigma = 0.5$	0.086	0.085	0.195	0.621
		$\sigma = 0.75$	0.160	0.160	0.340	1.431
		$\sigma = 1$	0.272	0.280	0.538	2.685
$n = 20$	$p = 101$	$\sigma = 0$	0.012	0.012	0.040	0.089
		$\sigma = 0.5$	0.047	0.046	0.117	0.660
		$\sigma = 0.75$	0.079	0.079	0.192	1.434
		$\sigma = 1$	0.135	0.138	0.301	2.549
$n = 50$	$p = 51$	$\sigma = 0$	0.025	0.025	0.070	0.076
		$\sigma = 0.5$	0.087	0.086	0.202	0.584
		$\sigma = 0.75$	0.162	0.169	0.348	1.435
		$\sigma = 1$	0.250	0.269	0.522	2.557
$n = 100$	$p = 101$	$\sigma = 0$	0.012	0.012	0.041	0.033
		$\sigma = 0.5$	0.045	0.045	0.114	0.518
		$\sigma = 0.75$	0.085	0.085	0.195	1.337
		$\sigma = 1$	0.138	0.139	0.299	2.618

Table 3: **Setting 3:** The MISE of the function estimates when K is large compared with n and p . The σ^2 -value is used to control the signal-to-noise ratio.

Sample Size	Dimension	Size of η	MISE			
			FASM	Bsmooth	Localin	PACE
$n = 50$	$p = 40$	$\sigma = 0$	0.097	0.097	0.121	0.225
		$\sigma = 0.5$	0.206	0.207	0.235	0.590
		$\sigma = 0.75$	0.291	0.309	0.335	1.044
		$\sigma = 1$	0.390	0.457	0.466	1.747

measurement errors.

We use B-spline basis functions of order four with knots at every data point. We show the residuals from the four models in Figure 5. Apart from the FASM, the other models produce large residuals, with some extreme values. To examine the residual structure from Bsmooth model, we conduct a principal component analysis on the residuals. Figure 4b shows the screeplot of the

eigenvalues. The residuals from Bsmooth model exhibit a “spike” structure, where the first few eigenvalues are significantly larger than the rest. This means the residuals contain information that can be captured by only a few factors, which calls for a further dimension reduction model on the residuals.

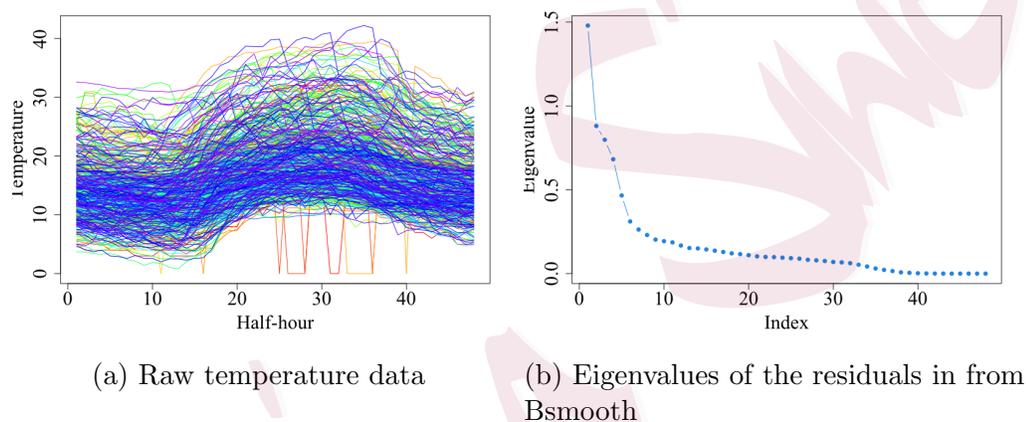


Figure 4: Half-hourly temperature data on Fridays at Adelaide airport.

7. Conclusion

In this paper, we propose an FASM for functional data. We study raw functional data, a mixture of functional curves, and high-dimensional errors. When the measurement error is informative, a smoothing model alone is inadequate for capturing data variation and recovering the signal functional component. The proposed model incorporates a factor structure into the smoothing model to further explain the large residuals. We propose a numerical iteration approach to obtain estimates in the smoothing and the factor models simultaneously. The asymptotic distribution of the estimators is given, with a proof. Our model also serves as a dimension-reduction method

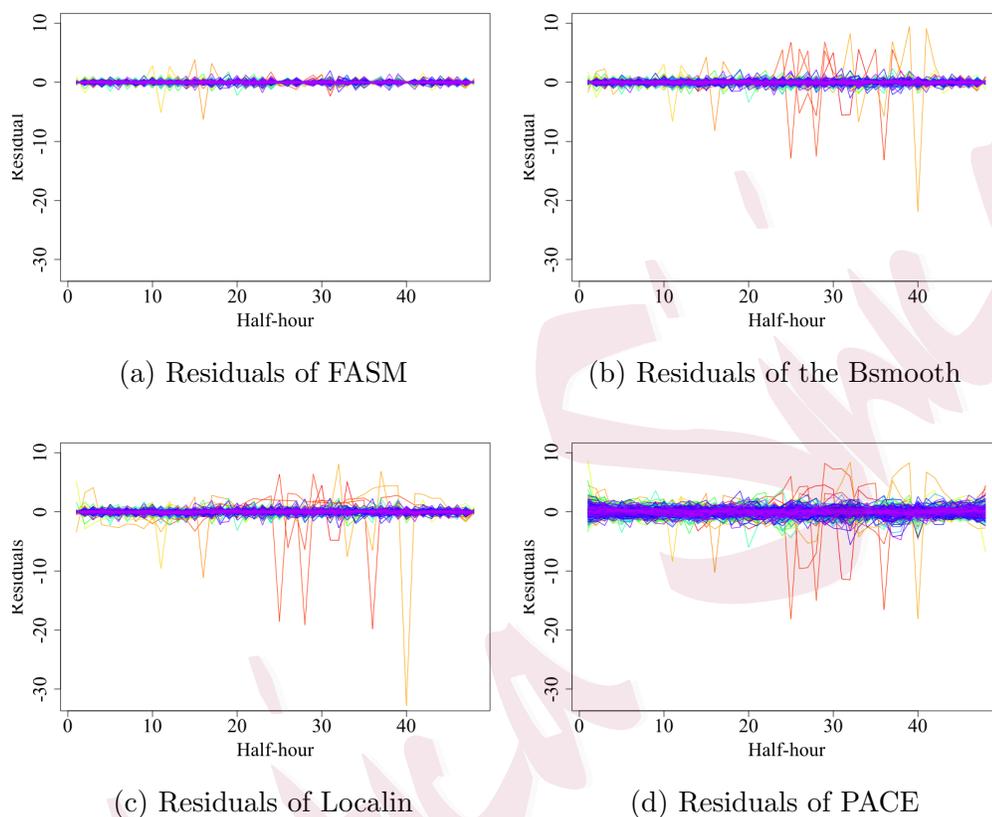


Figure 5: Half-hourly temperature data on Friday at Adelaide airport.

for functional and high-dimensional mixture data, easing the path to making inferences. We provide an example in which we construct a covariance estimator for the raw data. Furthermore, we show that the model can be applied when there is a misidentification in the data structure, two examples of which are the misspecification of the smoothing basis functions and neglecting the step jumps in the mean level of the functions. The advantages of the proposed model are demonstrated using simulation studies. We also show how our model performs by applying it to Australian temperature data.

The proposed model is a good starting point for modeling complex data structures. We deal

with a mixture of smooth functional curves and high-dimensional measurement errors. The factor model component can be regarded as a “boosting component” that improves model accuracy. Extending this idea, the model can be applied to other data structures, such data that contain change points. The change point is a popular problem in many statistics and econometric topics, and has been studied extensively in multivariate settings. Previous works on change points in functional data include those of [Berkes et al. \(2009\)](#), [Hörmann and Kokoszka \(2010\)](#), and [Hörmann et al. \(2015\)](#). In our simulations, we show examples that our model can be used to model functional data with change points in the cross-sectional direction. The model can be modified to account for change points in the sample direction. This is left to future research.

Supplementary Material

In the online Supplementary Material, we extend our model to nonparametric settings, introduce a covariance estimator, and present extensive simulated and real-data analysis. Here, we also provide proofs for the theorems and additional results.

References

- Ahn, S. C. and A. R. Horenstein (2013). Eigenvalue ratio test for the number of factors. *Econometrica* 81(3), 1203–1227.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6), 716–723.
- Bai, J. and S. Ng (2002). Determining the number of factors in approximate factor models. *Econometrica* 70(1), 191–221.
- Berkes, I., R. Gabrys, L. Horváth, and P. Kokoszka (2009). Detecting changes in the mean of functional observations.

- Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 71(5), 927–946.
- Cai, T. T. and M. Yuan (2011). Optimal estimation of the mean function based on discretely sampled functional data: Phase transition. *The Annals of Statistics* 39(5), 2330–2355.
- Cuevas, A. (2014). A partial overview of the theory of statistics with functional data. *Journal of Statistical Planning and Inference* 147, 1–23.
- Eubank, R. L. (1999). *Nonparametric Regression and Spline Smoothing* (2nd ed.). New York: Marcel Dekker.
- Fan, J., Y. Fan, and J. Lv (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics* 147(1), 186–197.
- Fan, J. and I. Gijbels (1996). *Local Polynomial Modelling and Its Applications*. London: Chapman & Hall.
- Febrero-Bande, M., P. Galeano, and W. González-Manteiga (2017). Functional principal component regression and functional partial least-squares regression: An overview and a comparative study. *International Statistical Review* 85(1), 61–83.
- Ferraty, F. and P. Vieu (2006). *Nonparametric Functional Data Analysis: Theory and Practice*. New York: Springer Science & Business Media.
- Goia, A. and P. Vieu (2016). An introduction to recent advances in high/infinite dimensional statistics. *Journal of Multivariate Analysis* 146, 1–6.
- Golub, G. H., M. Heath, and G. Wahba (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics* 21(2), 215–223.
- Green, P. J. and B. W. Silverman (1999). *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach*. London: Chapman & Hall.
- Hörmann, S., L. Kidziński, and M. Hallin (2015). Dynamic functional principal components. *Journal of the Royal Statistical*

- Society, Statistical Methodology, Series B* 77(2), 319–348.
- Hörmann, S. and P. Kokoszka (2010). Weakly dependent functional data. *The Annals of Statistics* 38(3), 1845–1884.
- Horváth, L. and P. Kokoszka (2012). *Inference for functional data with applications*, Volume 200. New York: Springer Science & Business Media.
- Jiang, B., Y. Yang, J. Gao, and C. Hsiao (2021). Recursive estimation in large panel data models: Theory and practice. *Journal of Econometrics* 224(2), 439–465.
- Knight, K. and W. Fu (2000). Asymptotics for lasso-type estimators. *The Annals of Statistics* 28(5), 1356–1378.
- Lam, C., Q. Yao, and N. Bathia (2011). Estimation of latent factors for high-dimensional time series. *Biometrika* 98(4), 901–918.
- Onatski, A. (2010). Determining the number of factors from empirical distribution of eigenvalues. *The Review of Economics and Statistics* 92(4), 1004–1016.
- Onatski, A. (2012). Asymptotics of the principal components estimator of large factor models with weakly influential factors. *Journal of Econometrics* 168(2), 244–258.
- Ramsay, J. O. and G. Hooker (2017). *Dynamic Data Analysis: Modeling Data with Differential Equations*. New York: Springer.
- Ramsay, J. O. and B. W. Silverman (2002). *Applied Functional Data Analysis*. New York: Springer.
- Ramsay, J. O. and B. W. Silverman (2005). *Functional Data Analysis*. New York: Springer.
- Reiss, P. T., J. Goldsmith, H. L. Shang, and R. T. Ogden (2017). Methods for scalar-on-function regression. *International Statistical Review* 85(2), 228–249.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics* 6(2), 461–464.
- Wahba, G. (1990). *Spline models for observational data*, Volume 59. SIAM.

- Wahba, G. and S. Ward (1975). Periodic splines for spectral density estimation: The use of cross validation for determining the degree of smoothing. *Communications in Statistics: Theory and methods* 4(2), 125–141.
- Wand, M. P. and C. M. Jones (1995). *Kernel Smoothing*. Boca Raton, FL: Chapman & Hall.
- Wang, J.-L., J.-M. Chiou, and H.-G. Müller (2016). Functional data analysis. *Annual Review of Statistics and Its Application* 3, 257–295.
- Yao, F., H. Müller, and J. Wang (2005). Functional data analysis for sparse longitudinal data. *Journal of the American Statistical Association* 100(470), 577–590.
- Yao, W. and R. Li (2013). New local estimation procedure for a non-parametric regression function for longitudinal data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 75(1), 123–138.
- Zhang, X. and J. L. Wang (2016). From sparse to dense functional data and beyond. *The Annals of Statistics* 44(5), 2281–2321.

Yuan Gao

E-mail: (yuan.gao@anu.edu.au)

Phone: (+61 2 612 57290)

Han Lin Shang

E-mail: (hanlin.shang@mq.edu.au)

Yanrong Yang

E-mail: (yanrong.yang@anu.edu.au)