

# Reconstructing partially observed functional data via factor models of increasing rank

Maximilian Ofner

Siegfried Hörmann

Graz University of Technology  
Institute of Statistics

COMPSTAT 2023  
London

- 1 Introduction
- 2 Estimation of reconstructions
- 3 Simultaneous prediction bands
- 4 Real data illustration

# Section 1

## Introduction

Let

$$X = (X(u) : u \in [0, 1])$$

be a (centred) random function *observable* on  $O \subset [0, 1]$ .

Let

$$X = (X(u) : u \in [0, 1])$$

be a (centred) random function *observable* on  $O \subset [0, 1]$ .

### Assumption (MCAR)

The set  $O$  is independent of  $X$ .

See [Liebl and Rameseder \(2019\)](#) for a relaxation of this assumption.

Let

$$X = (X(u) : u \in [0, 1])$$

be a (centred) random function *observable* on  $O \subset [0, 1]$ .

### Assumption (MCAR)

The set  $O$  is independent of  $X$ .

See [Liebl and Rameseder \(2019\)](#) for a relaxation of this assumption.

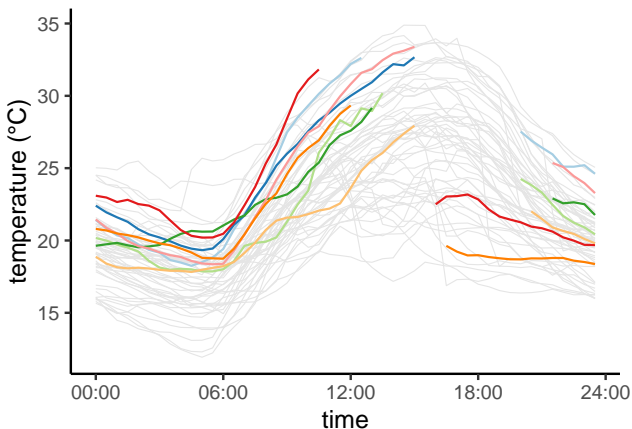
Assume  $X^O = (X(u) : u \in O)$  satisfies the Karhunen-Loève expansion

$$X^O(u) = \sum_{k=1}^{\infty} \xi_k^O \phi_k^O(u), \quad u \in O,$$

where

- $(\xi_k^O)$  are uncorrelated rv's with  $\mathbb{V}\text{ar}(\xi_k^O) = \lambda_k^O$ ,
- $(\phi_k^O)$  are ONB of  $L^2(O)$ .

Assume  $\lambda_1^O \geq \lambda_2^O \geq \dots$



**Figure:** Intraday temperature values measured half-hourly in the east of Graz (Austria) between 1 July and 31 August 2022. Complete observations in grey.

Data: [Land Steiermark \(2023\)](#)

## Question

How can we infer  $X$  from  $X^\circ$ ?



## Question

How can we infer  $X$  from  $X^\circ$ ?

Kneip and Liebl (2020) study the linear operator  $\mathcal{L} : L^2(\mathcal{O}) \rightarrow L^2[0, 1]$ ,

$$\mathcal{L}(X^\circ)(u) = \sum_{k=1}^{\infty} \xi_k^\circ \tilde{\phi}_k^\circ(u), \quad u \in [0, 1],$$

where

$$\tilde{\phi}_k^\circ(u) = \frac{\mathbb{E}[X(u)\xi_k^\circ]}{\lambda_k^\circ}, \quad u \in [0, 1],$$

are *extrapolated basis functions*.

## Question

How can we infer  $X$  from  $X^\circ$ ?

Kneip and Liebl (2020) study the linear operator  $\mathcal{L} : L^2(\mathcal{O}) \rightarrow L^2[0, 1]$ ,

$$\mathcal{L}(X^\circ)(u) = \sum_{k=1}^{\infty} \xi_k^\circ \tilde{\phi}_k^\circ(u), \quad u \in [0, 1],$$

where

$$\tilde{\phi}_k^\circ(u) = \frac{\mathbb{E}[X(u)\xi_k^\circ]}{\lambda_k^\circ}, \quad u \in [0, 1],$$

are *extrapolated basis functions*. It holds

$$X(u) = \mathcal{L}(X^\circ)(u) + Z^\circ(u), \quad u \in [0, 1],$$

for some *reconstruction error*  $Z^\circ = (Z^\circ(u) : u \in [0, 1])$ .

We want to avoid restrictive smoothness conditions such as differentiable paths.

## Idea

Estimate  $\mathcal{L}(X^O)$  via approximate factor models.

## Assumption (Rank)

There exists some  $r_O < \infty$  such that

$$X^O(u) = \sum_{k=1}^{r_O} \xi_k^O \varphi_k^O(u), \quad u \in O.$$

Later:  $r_O \rightarrow \infty$  to account for high dimensionality.

Let  $\{X_t : t \leq T\}$  be iid copies of  $X$ . Consider a regular grid (size  $N$ )

$$0 = u_1 < u_2 < \cdots < u_N = 1.$$

Let  $\{X_t : t \leq T\}$  be iid copies of  $X$ . Consider a regular grid (size  $N$ )

$$0 = u_1 < u_2 < \cdots < u_N = 1.$$

For  $e_{ti}$  independent of  $\{X_t : t \leq T\}$ , not necessarily iid, we observe

$$Y_{ti} = X_t(u_i) + e_{ti},$$

whenever  $u_i \in O_t \supset O$ .

Let  $\{X_t : t \leq T\}$  be iid copies of  $X$ . Consider a regular grid (size  $N$ )

$$0 = u_1 < u_2 < \dots < u_N = 1.$$

For  $e_{ti}$  independent of  $\{X_t : t \leq T\}$ , not necessarily iid, we observe

$$Y_{ti} = X_t(u_i) + e_{ti},$$

whenever  $u_i \in O_t \supset O$ . For  $u_i \in O$ , the rank assumption implies

$$Y_{ti} = \sum_{k=1}^{r_o} \xi_{tk}^o \phi_k^o(u_i) + e_{ti} = F_t^o \Lambda_i^o + e_{ti},$$

where we define

- $F_t^o = \left( \frac{\xi_{t1}^o}{\sqrt{\lambda_1^o}}, \dots, \frac{\xi_{tr_o}^o}{\sqrt{\lambda_{r_o}^o}} \right),$
- $\Lambda_i^o = \left( \sqrt{\lambda_1^o} \phi_1^o(u_i), \dots, \sqrt{\lambda_{r_o}^o} \phi_{r_o}^o(u_i) \right)'.$

# Approximate factor model

If the  $e_{ti}$  are only mildly correlated, then

$$Y_{ti} = F_t^\circ \Lambda_i^\circ + e_{ti},$$

constitutes an *approximate factor model* (AFM) of rank  $r_0$ .

A major branch of works analyses AFMs under a double asymptotic

$$N, T \rightarrow \infty.$$

Important contributions include [Chamberlain and Rothschild \(1983\)](#), [Bai and Ng \(2002\)](#), [Bai \(2003\)](#), and [Fan et al. \(2013\)](#). Imputation of missing values is studied by [Bai and Ng \(2021\)](#), [Cahan et al. \(2023\)](#), and [Xiong and Pelger \(2023\)](#).

# Important features of our factor model

- The rank  $r_0$  depends on  $O \in [0, 1]$ .
- We allow the rank  $r_0$  to grow and consider a triple asymptotic

$$r_0, N, T \rightarrow \infty.$$

See [Li et al. \(2017\)](#) and [Hörmann and Jammoul \(2022\)](#) for a similar asymptotic.

- The eigenvalue  $\lambda_k^O$  measures the pervasiveness of the  $k$ -th factor score and enters our convergence rates.

See [Bai and Ng \(2023\)](#) for related results on AFMs with weaker loadings.



## Section 2

# Estimation of reconstructions

- 1 Let  $\widehat{F}_t^o$  ( $1 \times r_o$ ) be the PC-estimate of the  $t$ -th factors.
- 2 Assume the first  $T_c \leq T$  curves are completely observable.
  - $\widehat{F}^o \dots (T_c \times r_o)$  matrix with  $\widehat{F}_t^o$  in its rows,  $t \leq T_c$ .
  - $Y_i \dots (1 \times T_c)$  vector with elements  $Y_{ti}$ ,  $t \leq T_c$ .

Convention:  $\frac{\widehat{F}^{o'} \widehat{F}^o}{T_c} = I$ .

- 3 Estimate  $\mathcal{L}(X_t^o)(u_i)$  by projecting  $Y_i$  onto the estimated factors,

$$\widehat{\mathcal{L}(X_t^o)}(u_i) = \widehat{F}_t^o (\widehat{F}^{o'} \widehat{F}^o)^{-1} \widehat{F}^{o'} Y_i = \frac{\widehat{F}_t^o \widehat{F}^{o'} Y_i}{T_c}$$

and interpolate linearly.

## Theorem (OH, 2023)

Let  $X$  be observable on  $O \subset [0, 1]$ . Under regularity conditions,

$$\sup_{u \in [0, 1]} |\widehat{\mathcal{L}(X^O)}(u) - \mathcal{L}(X^O)(u)| = O_p \left( \frac{r_O}{\lambda_{r_O}^O} \sqrt{\frac{1}{N_O} + \frac{\log(N)}{T_C}} \right),$$

as  $r_O \rightarrow \infty$ ,  $N_O \rightarrow \infty$  with  $N \rightarrow \infty$  and  $T_C \rightarrow \infty$  with  $T \rightarrow \infty$ .

- $r_O$  ... rank of  $X^O$  on  $O$ ,
- $N_O$  ... number of grid points  $u_i \in O$ ,
- $T_C$  ... number of complete curves.

If  $O = [0, 1]$ , then  $\mathcal{L}(X^O) = X$  and the estimator  $\hat{X} = \widehat{\mathcal{L}(X^O)}$  can be seen as a smoother of  $X$ .

## Theorem (OH, 2023)

Let  $X$  be observable on  $O \subset [0, 1]$ . Under regularity conditions,

$$\sup_{u \in [0, 1]} |\widehat{\mathcal{L}(X^O)}(u) - \mathcal{L}(X^O)(u)| = O_p \left( \frac{r_O}{\lambda_{r_O}^O} \sqrt{\frac{1}{N_O} + \frac{\log(N)}{T_C}} \right),$$

as  $r_O \rightarrow \infty$ ,  $N_O \rightarrow \infty$  with  $N \rightarrow \infty$  and  $T_C \rightarrow \infty$  with  $T \rightarrow \infty$ .

$r_O$  ... rank of  $X^O$  on  $O$ ,

$N_O$  ... number of grid points  $u_i \in O$ ,

$T_C$  ... number of complete curves.

If  $O = [0, 1]$ , then  $\mathcal{L}(X^O) = X$  and the estimator  $\hat{X} = \widehat{\mathcal{L}(X^O)}$  can be seen as a smoother of  $X$ .

## Theorem (OH, 2023)

Let  $X$  be observable on  $O \subset [0, 1]$ . Under regularity conditions,

$$\sup_{u \in [0, 1]} |\widehat{\mathcal{L}(X^O)}(u) - \mathcal{L}(X^O)(u)| = O_p \left( \frac{r_O}{\lambda_{r_O}^O} \sqrt{\frac{1}{N_O} + \frac{\log(N)}{T_C}} \right),$$

as  $r_O \rightarrow \infty$ ,  $N_O \rightarrow \infty$  with  $N \rightarrow \infty$  and  $T_C \rightarrow \infty$  with  $T \rightarrow \infty$ .

- $r_O$  ... rank of  $X^O$  on  $O$ ,
- $N_O$  ... number of grid points  $u_i \in O$ ,
- $T_C$  ... number of complete curves.

If  $O = [0, 1]$ , then  $\mathcal{L}(X^O) = X$  and the estimator  $\hat{X} = \widehat{\mathcal{L}(X^O)}$  can be seen as a smoother of  $X$ .

## Theorem (OH, 2023)

Let  $X$  be observable on  $O \subset [0, 1]$ . Under regularity conditions,

$$\sup_{u \in [0, 1]} |\widehat{\mathcal{L}(X^O)}(u) - \mathcal{L}(X^O)(u)| = O_p \left( \frac{r_O}{\lambda_{r_O}^O} \sqrt{\frac{1}{N_O} + \frac{\log(N)}{T_C}} \right),$$

as  $r_O \rightarrow \infty$ ,  $N_O \rightarrow \infty$  with  $N \rightarrow \infty$  and  $T_C \rightarrow \infty$  with  $T \rightarrow \infty$ .

- $r_O$  ... rank of  $X^O$  on  $O$ ,
- $N_O$  ... number of grid points  $u_i \in O$ ,
- $T_C$  ... number of complete curves.

If  $O = [0, 1]$ , then  $\mathcal{L}(X^O) = X$  and the estimator  $\hat{X} = \widehat{\mathcal{L}(X^O)}$  can be seen as a smoother of  $X$ .

## Section 3

# Simultaneous prediction bands

## Question

How can we construct simultaneous prediction bands (SPB) for the reconstructions?

Assume  $\widehat{X}_t$  and  $\widehat{\mathcal{L}(X_t^0)}$  consistently estimate  $X_t$  and  $\mathcal{L}(X_t^0)$ . Set

$$\widehat{Z}_t^0(u) = \begin{cases} \widehat{X}_t(u) - \widehat{\mathcal{L}(X_t^0)}(u), & u \in M = [0, 1] \setminus O, \\ 0, & u \in O, \end{cases}$$

and let  $\widehat{q}_\alpha$  be an estimator of  $q_\alpha$  defined by

$$\mathbb{P}(\sup_{u \in M} \{|Z^0(u)|/\text{sd}(Z^0(u))\} > q_\alpha) = \alpha.$$

For  $\alpha \in (0, 1)$ , we then consider the SPB

$$\widehat{\mathcal{L}(X_t^0)}(u) \pm \widehat{q}_\alpha \widehat{\text{sd}}(Z^0(u)), \quad u \in M.$$



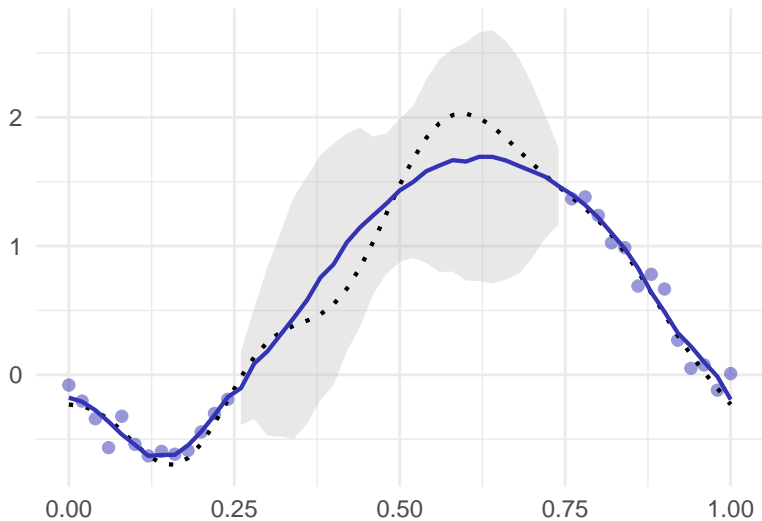
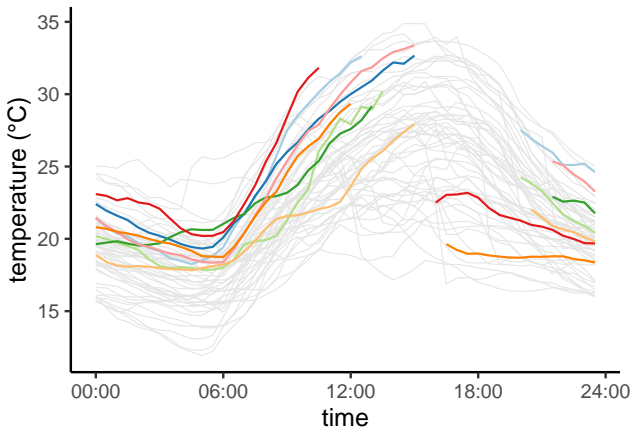


Figure: Reconstruction of a simulated curve along with 95% prediction band.

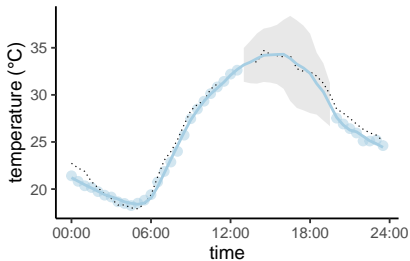
## Section 4

### Real data illustration

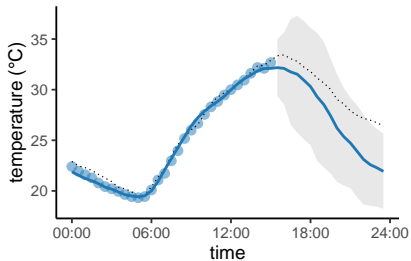


**Figure:** Intraday temperature values measured half-hourly in the east of Graz (Austria) between 1 July and 31 August 2022. Complete observations in grey.

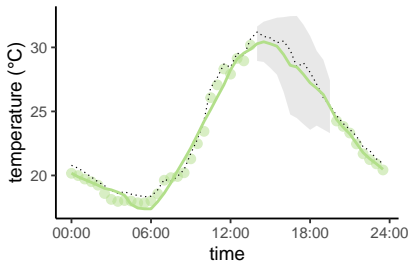
2022-07-21



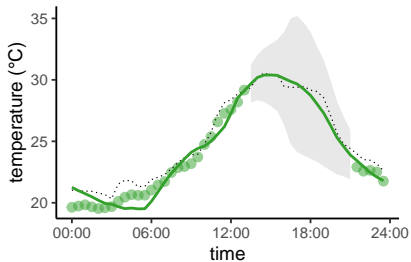
2022-07-25



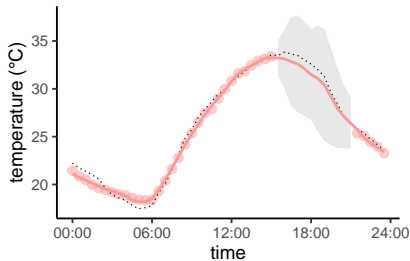
2022-08-01



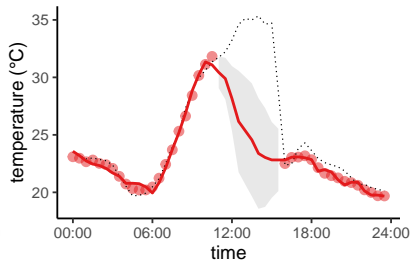
2022-08-16



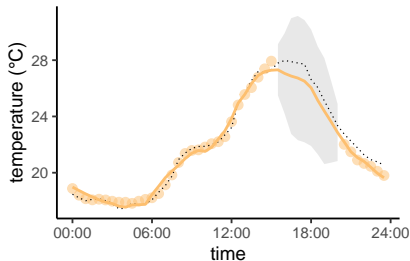
2022-08-17



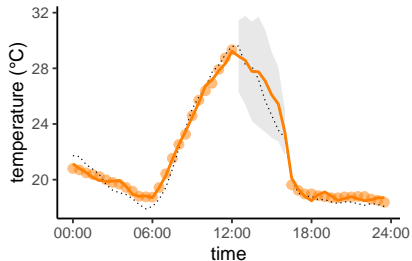
2022-08-18



2022-08-25



2022-08-27



# Questions? Comments? Criticism?

Manuscript:

► [arXiv:2305.13152](https://arxiv.org/abs/2305.13152)

Code:

► [Github: FDFM](#)

Maximilian Ofner  
[m.ofner@tugraz.at](mailto:m.ofner@tugraz.at)



# References

- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica*, 71(1):135–171.
- Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.
- Bai, J. and Ng, S. (2021). Matrix completion, counterfactuals, and factor analysis of missing data. *J. Amer. Statist. Assoc.*, 116(536):1746–1763.
- Bai, J. and Ng, S. (2023). Approximate factor models with weaker loadings. *Journal of Econometrics*.
- Cahan, E., Bai, J., and Ng, S. (2023). Factor-based imputation of missing values and covariances in panel data of large dimensions. *J. Econometrics*, 233(1):113–131.
- Chamberlain, G. and Rothschild, M. (1983). Arbitrage, factor structure, and mean-variance analysis on large asset markets. *Econometrica*, 51(5):1281–1304.
- Fan, J., Liao, Y., and Mincheva, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 75(4):603–680.
- Hörmann, S. and Jammoul, F. (2022). Consistently recovering the signal from noisy functional data. *J. Multivariate Anal.*, 189:Paper No. 104886, 18.
- Kneip, A. and Liebl, D. (2020). On the optimal reconstruction of partially observed functional data. *Ann. Statist.*, 48(3):1692–1717.
- Land Steiermark (2023). Air quality data. <https://app.luis.steiermark.at/luft2/suche.php>. Licensed under CC BY 4.0; accessed on March 28, 2023.
- Li, H., Li, Q., and Shi, Y. (2017). Determining the number of factors when the number of factors can increase with sample size. *J. Econometrics*, 197(1):76–86.
- Liebl, D. and Rameseder, S. (2019). Partially observed functional data: the case of systematically missing parts. *Comput. Statist. Data Anal.*, 131:104–115.
- Xiong, R. and Pelger, M. (2023). Large dimensional latent factor modeling with missing observations and applications to causal inference. *J. Econometrics*, 233(1):271–301.
- Yao, F., Müller, H.-G., and Wang, J.-L. (2005). Functional data analysis for sparse longitudinal data. *J. Amer. Statist. Assoc.*, 100(470):577–590.