Registration for Incomplete Non-Gaussian Functional Data

Dr. Alexander Bauer LMU Munich, Germany









Fabian Scheipl

Helmut Küchenhoff



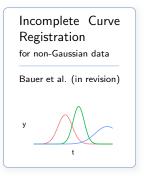
Alice-Agnes Gabriel

University of Colorado





Julia Wrobel





- 1. Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package

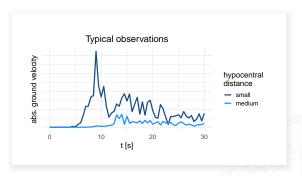
Incomplete Curve Registration for non-Gaussian data Bauer et al. (in revision)

registr 2.0 Wrobel & Bauer (2021)

1. Conceptual Basics

- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- Simulation Study
- 4. Implementation in registr package

Application with data on seismic ground motion propagation



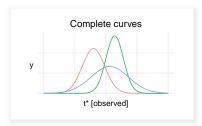
Research Question

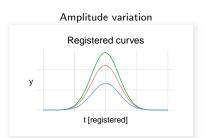
Data

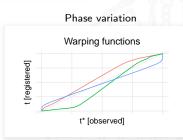
Given the occurrence of a seismic event, what are the driving forces for its strength?

135 simulations of the 1994 Northridge (US) quake, 30s recordings of ground motion at \sim 6000 seismometers

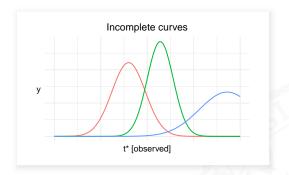
Separating amplitude and phase



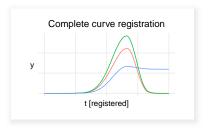


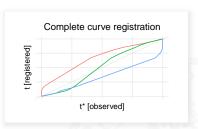


Challenge incomplete curves



Challenge incomplete curves





Common assumption: Processes are observed until their very end

We want a method that ...

- is able to handle incomplete curves
- is able to handle non-Gaussian data
- includes a lower-dimensional representation of the registered curves

Challenge 1

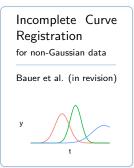
Lack of registration methods for incomplete curves

Challenge 2

Registration methods often only applicable to specific data situations Complete, Gaussian, densely observed curves on a regular grid

Challenge 3

Lack of good software packages for incomplete curve registration.



registr 2.0 Wrobel & Bauer (2021)

- 1. Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package

Incomplete Approach

Notation

$$Y_{i}(t_{i}^{*})$$
 $T = [0, 1]$
 $h_{i}^{-1}: T_{i}^{*} \mapsto T$
 $Y_{i}(t) = Y_{i}(h_{i}^{-1}(t_{i}^{*}))$

curve of subject i = 1, ..., N, observed over chronological time domain $T_i^* = [t_{\min,i}^*, t_{\max,i}^*]$ *internal time domain*, with $T_i^* \subseteq T \ \forall i$ $h_i^{-1}: T_i^* \mapsto T$ inverse warping functions $Y_i(t) = Y_i(h_i^{-1}(t_i^*))$ aligned curve of subject i

Alexander Bauer 5 / 20

Notation

$$Y_i(t_i^*)$$
 $T = [0, 1]$
 $h_i^{-1} : T_i^* \mapsto T$
 $Y_i(t) = Y_i(h_i^{-1}(t_i^*))$

n;

$$\begin{split} T_i^* &= [t_{\min,i}^*, 1] \\ T_i^* &= [0, t_{\max,i}^*] \\ T_i^* &= [t_{\min,i}^*, t_{\max,i}^*] \end{split}$$

curve of subject $i=1,\ldots,N$, observed over chronological time domain $T_i^*=[t_{\min,i}^*,t_{\max,i}^*]$ internal time domain, with $T_i^*\subseteq T\ \forall\ i$ inverse warping functions aligned curve of subject i

number of measurements for curve i

leading incompleteness trailing incompleteness full incompleteness

Incomplete Approach

Only few approaches exist for registering incomplete curves

- Dynamic Time Warping methods developed for matching time series
 - ? Computationally quite expensive & mostly very heuristic algorithms
- Bayesian approach for fragmented functional data Matuk et al. (2019)
 Computationally quite expensive
- SRVF-based approach for elastic partial matching Bryner & Srivastava (2021)
 based on the Square-Root Velocity Function framework
 - o Joint estimation scheme for complete curve SRVF and linear domain scaling
 - £ Limited to continuous data & based on derivatives

Incomplete Approach

Only few approaches exist for registering incomplete curves

- Dynamic Time Warping methods developed for matching time series
 - f Computationally quite expensive & mostly very heuristic algorithms
- Bayesian approach for fragmented functional data Matuk et al. (2019)
 - Computationally quite expensive
- SRVF-based approach for elastic partial matching Bryner & Srivastava (2021)
 based on the Square-Root Velocity Function framework
 - Joint estimation scheme for complete curve SRVF and linear domain scaling
 - £ Limited to continuous data & based on derivatives

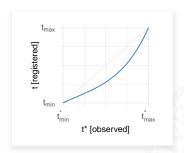
Extend the complete curve likelihood-based framework of Wrobel et al. (2019)

- ✓ Applicable for exponential family distributions
- ✓ Joint estimation scheme with Gaussian or Binomial FPCA

Aim: Map the chronological domains T_i^* onto the registered domain T

 \Rightarrow Inverse warping functions h_i^{-1} map curve $Y_i(t_i^*)$ onto a template $\mu_i(t)$:

$$\mathsf{E}\left[Y_{i}\left(h_{i}^{-1}(t_{i}^{*})\right)|h_{i}^{-1}\right]=\mu_{i}\left(t\right).$$



⇒ Warping functions have to be monotone and domain-preserving

Incomplete Approach

Low-dimensional **B-spline representation** of inverse warping functions:

$$h_i^{-1}(t_i^*) = \boldsymbol{\Theta}_h(t_i^*)\boldsymbol{\beta}_i$$

Notation

 Θ_{h} design matrix of K_{h} basis functions, $\in \mathbb{R}_{n_{i} \times K_{h}}$

 β_i coefficient vector, $\in \mathbb{R}_{K_b \times 1}$

Low-dimensional **B-spline representation** of inverse warping functions:

$$h_i^{-1}(t_i^*) = \boldsymbol{\Theta}_h(t_i^*)\boldsymbol{\beta}_i$$

Optimization of the **exponential family log-likelihood** for curve i:

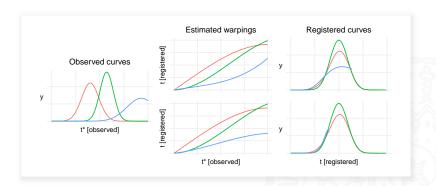
$$\ell\left(h_i^{-1}|y_i,\mu_i
ight) = \log\left(\prod_{j=1}^{n_i} f_{i,j}\left[y_i(t_{i,j}^*)
ight]
ight),$$

with $f_{i,j}(\cdot)$ the corresponding density with expected value $\mu_i\left(h_i^{-1}(t_{i,j}^*)\right)$, and conditional independence

- across functions $[Y_i \perp Y_{i'}] | \mu_i, \mu_{i'},$
- within functions $[Y_i(t_{ii}) \perp Y_i(t_{ik})] | \mu_i$.

Ensure reasonable warping functions $h_i^{-1}(t_i^*)$

- A constrained optimizer ensures monotony and domain-preservation ✓
- Further, they should produce scientifically reasonable distortions



Circumventing fixed time intervals

- \bullet Allow warping functions to start and / or end anywhere in the overall domain
- Penalize how much the duration of the (observed) time domain is warped

Penalized log-likelihood for full incompleteness

$$\begin{split} \ell_{\mathrm{pen}}\left(h_i^{-1}|y_i,\mu_i\right) &= \ell\left(h_i^{-1}|y_i,\mu_i\right) - \lambda \cdot n_i \cdot \mathrm{pen}\left(h_i^{-1}\right),\\ \text{with} \qquad & \mathrm{pen}\left(h_i^{-1}\right) = \left(\left[h_i^{-1}(t_{\mathrm{max},i}^*) - h_i^{-1}(t_{\mathrm{min},i}^*)\right] - \left[t_{\mathrm{max},i}^* - t_{\mathrm{min},i}^*\right]\right)^2. \end{split}$$

Circumventing fixed time intervals

- \bullet Allow warping functions to start and / or end anywhere in the overall domain
- Penalize how much the duration of the (observed) time domain is warped

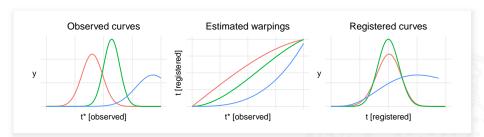
Penalized log-likelihood for full incompleteness

$$\begin{split} \ell_{\mathrm{pen}}\left(h_i^{-1}|y_i,\mu_i\right) &= \ell\left(h_i^{-1}|y_i,\mu_i\right) - \lambda \cdot n_i \cdot \mathrm{pen}\left(h_i^{-1}\right),\\ \text{with} \qquad & \mathrm{pen}\left(h_i^{-1}\right) = \left(\left[h_i^{-1}(t_{\mathrm{max},i}^*) - h_i^{-1}(t_{\mathrm{min},i}^*)\right] - \left[t_{\mathrm{max},i}^* - t_{\mathrm{min},i}^*\right]\right)^2. \end{split}$$

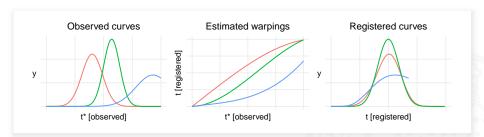
Simplification for trailing incompleteness (with $h_i^{-1}(t_{\min,i}^*) = t_{\min,i}^* \, \forall \, i$):

$$\mathsf{pen}\left(h_i^{-1}\right) = \left[h_i^{-1}(t_{\mathsf{max},i}^*) - t_{\mathsf{max},i}^*\right]^2.$$

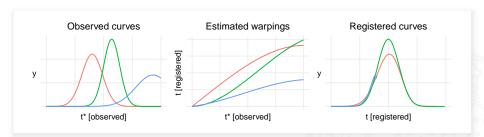
Example: Using a way too high value, $\lambda = 10$



Example: Using a too high value, $\lambda = 1$

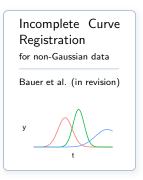


Example: Using a **too low value**, $\lambda = 0$



Example: Using a **fitting value**, $\lambda = 0.25$





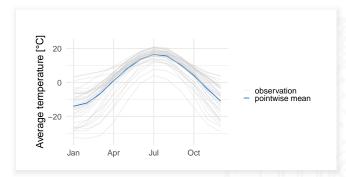
registr 2.0 Wrobel & Bauer (2021)

- 1. Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package

Concept of Functional Principal Component Analysis

- Estimation of main modes of variation around the mean curve
- Can be performed similarly to scalar PCA taking the eigenvalues of the covariance structure

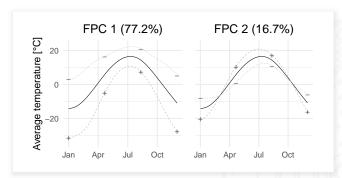
Example: Canadian weather data Ramsay & Silverman (2005)



Concept of Functional Principal Component Analysis

- Estimation of main modes of variation around the mean curve
- Can be performed similarly to scalar PCA taking the eigenvalues of the covariance structure

Example: Canadian weather data Ramsay & Silverman (2005)



Generalized FPCA

Adapt the two-step approach of Gertheiss et al. (2017)

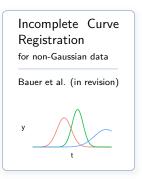
Combination of a nonparametric covariance estimator and a Functional Mixed Model

- ✓ Applicable to diverse exponential family settings
- ✓ Availability of efficient, robust software

Central sources of bias for incomplete curves

- Poor coverage of the overall domain
- Violation of MCAR assumption
- ⇒ (Severe) bias of mean and covariance estimators Liebl & Rameseder (2019)

▶ details on GFPCA estimation



registr 2.0 Wrobel & Bauer (2021)

- Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package

Iterative estimation algorithm Wrobel et al. (2019)

Aim:

- Registration of all observed curves $Y_i(t_i^*)$ to a comparable shape
- Low-rank representation of registered curves $Y_i(t) = Y_i(h_i^{-1}(t_i^*))$

Joint Approach

Iterative estimation algorithm Wrobel et al. (2019)

1. Initialize $\hat{h}_i^{-1}(t^*)^{[0]}$: Register curves $y_i(t_i^*)$ to initial template $\mu(t)^{[0]}$

Iterative estimation algorithm Wrobel et al. (2019)

- 1. Initialize $\hat{h}_i^{-1}(t^*)^{[0]}$: Register curves $y_i(t_i^*)$ to initial template $\mu(t)^{[0]}$
- 2. Iterate over index $q = 0, 1, \dots$ while $\sum_{i=1}^{N} \left(\sum_{j=1}^{n_i} \left[\hat{h}_i^{-1}(t_{i,j}^*)^{[q]} - \hat{h}_i^{-1}(t_{i,j}^*)^{[q-1]} \right]^2 \right) > \Delta_h$
 - (i) Update GFPCA using registered curves $y_i\left(\hat{h}_i^{-1}(m{t}_i^*)^{[q-1]}\right)$
 - (ii) Re-estimate GFPCA representations $\mu_i(t)^{[q]}$ based on first $K^{[q]}$ FPCs explaining a share κ_{var} of the total variance
 - (iii) Update estimates $\hat{h}_i^{-1}(t^*)^{[q]}$ by re-registering curves $y_i(t_i^*)$ to $\mu_i(t)^{[q]}$

▶ details on template function choice

Iterative estimation algorithm Wrobel et al. (2019)

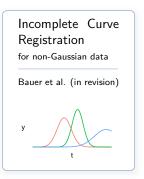
- 1. Initialize $\hat{h}_i^{-1}(t^*)^{[0]}$: Register curves $y_i(t_i^*)$ to initial template $\mu(t)^{[0]}$
- 2. Iterate over index $q = 0, 1, \dots$ while $\sum_{i=1}^{N} \left(\sum_{j=1}^{n_i} \left[\hat{h}_i^{-1}(t_{i,j}^*)^{[q]} - \hat{h}_i^{-1}(t_{i,j}^*)^{[q-1]} \right]^2 \right) > \Delta_h$

(i) Update GFPCA using registered curves
$$v_i \left(\hat{h}_i^{-1}(\mathbf{t}_i^*)^{[q-1]} \right)$$

- (i) Update GFPCA using registered curves $y_i \left(\hat{h}_i^{-1} (t_i^*)^{[q-1]} \right)$
- (ii) Re-estimate GFPCA representations $\mu_i(t)^{[q]}$ based on first $K^{[q]}$ FPCs explaining a share κ_{var} of the total variance
- (iii) Update estimates $\hat{h}_i^{-1}(t^*)^{[q]}$ by re-registering curves $v_i(t_i^*)$ to $\mu_i(t)^{[q]}$
- 3. Final GFPCA estimation based on registered curves $y_i\left(\hat{h}_i^{-1}(t_i^*)^{[q]}\right)$
 - \Rightarrow representations $\mu_i(t)$, based on K FPCs explaining a variance share $> \kappa_{\rm var}$

details on template function choice

Simulation Study

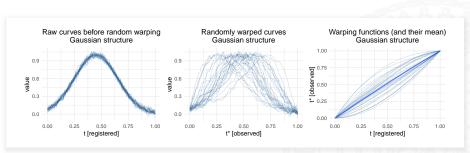




- 1. Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package

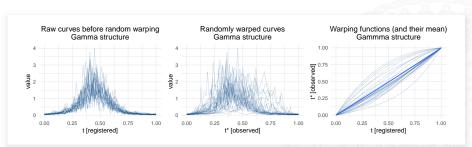
Simulation study

- Basis: Randomly warped (Gamma-transformed) Gaussian density each setting is run 100 times with 100 curves
- Simulated weak or strong trailing incompleteness Strong: Cut-off simulated in last 70% of the domain
- Amplitude rank \in {1, 2–3, 3–4} also because perfect identifiability is only given for amplitude variation of rank 1

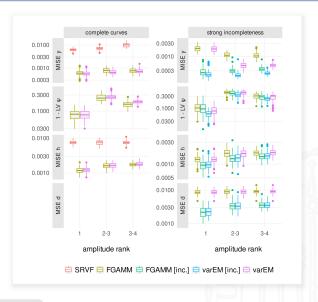


Simulation study

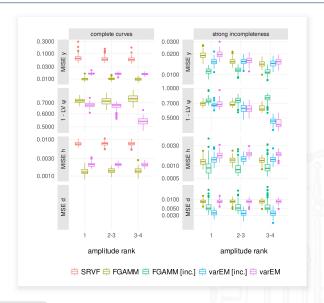
- Basis: Randomly warped (Gamma-transformed) Gaussian density each setting is run 100 times with 100 curves
- Simulated weak or strong trailing incompleteness Strong: Cut-off simulated in last 70% of the domain
- Amplitude rank \in {1, 2–3, 3–4} also because perfect identifiability is only given for amplitude variation of rank 1



Simulation Study – Gaussian structure



Simulation Study – Gamma structure



Simulation Study

Results

- Overall: Incomplete curve methods better represent the joint variation structure
- Phase:
 Incomplete curve methods better estimate the warping structure
- Amplitude:
 FGAMM struggles with the estimation of the FPC structure

Simulation Study

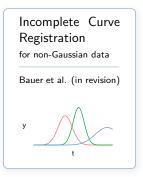
Results

- Overall: Incomplete curve methods better represent the joint variation structure
- Phase: Incomplete curve methods better estimate the warping structure
- Amplitude:
 FGAMM struggles with the estimation of the FPC structure

Further results

- Results are similar for weak and strong incompleteness
- Results are similar for settings with correlated amplitude and phase, and correlated amplitude and incompleteness
- FGAMM approach computationally quite inefficient
 Gamma runtime on 3000 curves, each with 50 measurements: ~ 0:27h

▶ details on runtimes

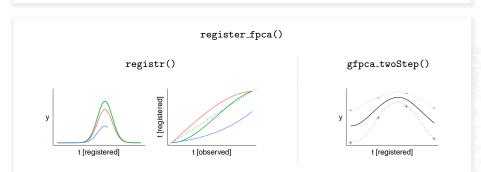


registr 2.0 Wrobel & Bauer (2021)

- 1. Conceptual Basics
- 2. Novel Approach
 - 2.1. Incomplete Curve Registration
 - 2.2. Incomplete Curve GFPCA
 - 2.3. Joint Approach
- 3. Simulation Study
- 4. Implementation in registr package



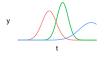
- Joint registration and GFPCA
- Applicable for leading / trailing / full incompleteness



Incomplete Curve Registration

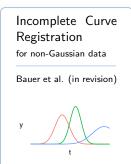
for non-Gaussian data

Bauer et al. (in revision)



registr 2.0
Wrobel & Bauer (2021)

- Novel approach for incomplete curves handling leading / trailing / full incompleteness
- Ability to handle non-Gaussian curves on irregular, potentially sparse grids
- √ registr package
 applicable to diverse data settings

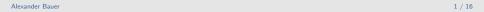


registr 2.0
Wrobel & Bauer (2021)

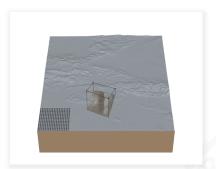
- Novel approach for incomplete curves handling leading / trailing / full incompleteness
- Ability to handle non-Gaussian curves on irregular, potentially sparse grids
- √ registr package
 applicable to diverse data settings

- Robust & intuitive penalization
- Robust & efficient covariance estimation
- --→ Analysis of seismic amplitude and phase

Appendix



Application with data on seismic ground motion propagation



Research Question

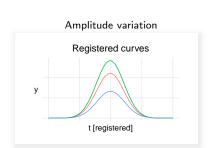
Data

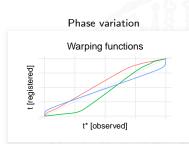
Given the occurrence of a seismic event, what are the driving forces for its strength?

135 simulations of the 1994 Northridge (US) quake, 30s recordings of ground motion at ~ 6000 seismometers

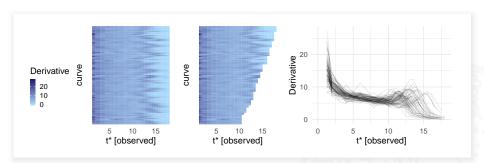
Basic structure of registration algorithms

- 1. Choose a template function
- 2. Choose a reasonable objective function
- 3. Optimize wrt. ensuring the well-definedness of warping functions

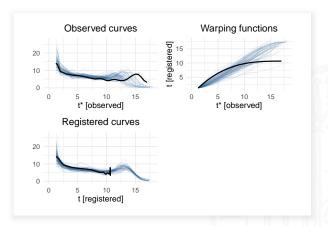




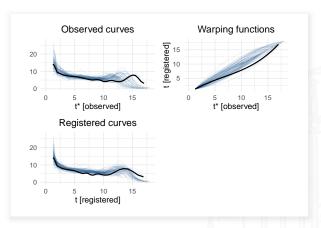
Example: Berkeley child growth data with simulated trailing incompleteness



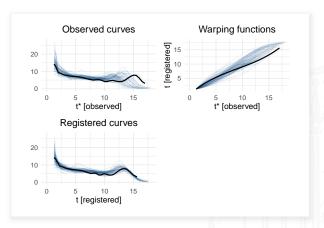
Example: **Berkeley child growth data** with simulated trailing incompleteness, using a **too small value**, $\lambda = 0$.



Example: Berkeley child growth data with simulated trailing incompleteness, using a too high value, $\lambda = 100$.



Example: Berkeley child growth data with simulated trailing incompleteness, using a reasonable value, $\lambda = 0.8$.



Two-step approach

applied to the registered curves $Y_i(t) = Y_i(h_i^{-1}(t_i^*))$

- 1. Estimation of FPCs $\psi_k(t)$ based on a marginal method Hall et al. (2008)
- 2. Estimation of mean $\alpha(t)$ and FPC scores c_i through a Generalized Functional Additive Mixed Model Gertheiss et al. (2017)

$$E[Y_i(t)] = \mu_i(t) = g[X_i(t)],$$

$$X_i(t) \approx \alpha(t) + \sum_{k=1}^K c_{i,k} \cdot \psi_k(t),$$

Notation

K number of principal components $\mu_i(t)$ conditional expected value of $Y_i(t)$

 $g[X_i(t)]$ latent Gaussian process transformed with response function $g(\cdot)$

details on GEPCA estimation

Step 1 – Estimation of FPCs Hall et al. (2008)

based on $E[Y_i(t)] = g[X_i(t)]$

- 1. Center curves $Y_i(t)$ based on a marginal estimate of $\mu_Y(t) = \mathbb{E}[Y_i(t)]$ by smoothing the data in a generalized additive model
- 2. Marginal estimation of the covariance:

$$\widehat{\mathsf{Cov}}\left[X_i(s),X_i(t)
ight]pprox rac{\hat{\sigma}_Y(s,t)}{g^{(1)}[\mu_X(s)]\cdot g^{(1)}[\mu_X(t)]},$$

with

- o $\sigma_Y(s,t) = \mathsf{E}[Y_{c,i}(s) \cdot Y_{c,i}(t)]$ based on centered curves $Y_{c,i}(t)$, with $\sigma_Y(s_1,s_2)$ the mean of all pairwise products $y_{c,i}(s_1) \cdot y_{c,i}(s_2)$, and $\hat{\sigma}_Y(s,t)$ a smoothed version of $\sigma_Y(s,t)$ using a tensor product P-spline basis
- the marginal mean $\mu_X(t)$ estimated accordingly to $\mu_Y(t)$,
- o $g^{(1)}(\cdot)$ the first derivative of the response function.
- 3. Spectral decomposition to yield FPCs $\psi_k(t)$ and associated eigenvalues τ_k

▶ back to GFPCA estimation

Step 2 - Estimation of FPC scores Gertheiss et al. (2017)

Estimation of mean $\alpha(t)$ and FPC scores c_i conditional on $\psi_k(t)$ in a **Generalized Functional Additive Mixed Model**:

$$g\left[\mu_i(t)\right] = \alpha(t) + \sum_{k=1}^K c_{i,k} \cdot \psi_k(t),$$

with the $c_{i,k} \sim N(0, \tau_k)$ random effects in an FPC basis representation.

⇒ Use robust routines (gamm4 / lme4), highly efficient for many random effects

Notation Notation

 $\alpha(t) = \Theta_{\alpha} \alpha$ smooth effect, with P-spline basis Θ_{α} and parameters α

▶ back to GFPCA estimation

Approaches for non-Gaussian FPCA

Most existing approaches either assume Gaussianity Stefanucci et al. (2018) or perform a marginal, potentially biased estimation Gertheiss et al. (2017)

Adapt the two-step approach of Gertheiss et al. (2017)

Combination of a nonparametric covariance estimator and a Functional Mixed Model

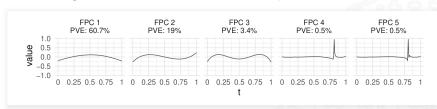
- ✓ Applicable to diverse exponential family settings
- ✓ Availability of efficient, robust software

Practical considerations

- Central sources of bias
 - Poor coverage of the overall domain
 - Violation of MCAR assumption
 - ⇒ (Severe) bias of mean and covariance estimators Liebl & Rameseder (2019)
- Choosing the number of FPCs based on explained variance share κ_{var}
 - o Explained shares of variance refer to the smoothed covariance surface
 - o Spectral decompositions often yield many subordinate FPCs

Practical considerations

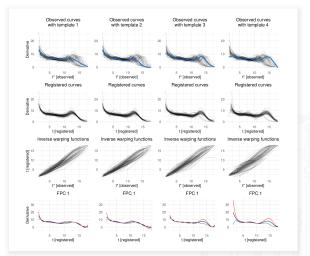
- Central sources of bias
 - Poor coverage of the overall domain
 - Violation of MCAR assumption
 - ⇒ (Severe) bias of mean and covariance estimators Liebl & Rameseder (2019)
- Choosing the number of FPCs based on explained variance share κ_{var}



⇒ Exclude such subordinate FPCs with minor explained shares of variance

Joint Approach

Choice of initial template function



back to the joint algorithm

Simulation Study

Compared methods all performing joint registration and FPCA

SRVF Complete curve approach of Tucker (2014)

FGAMM Complete curve approach based on two-step GFPCA

FGAMM [inc.] \hookrightarrow adapted for incomplete curves

varEM Complete curve approach of Wrobel et al. (2019)

varEM [inc.] \hookrightarrow adapted for incomplete curves

Performance metrics based on Mean (Integrated) Squared Errors

MISE_y Comparison of the simulated mean curves (before adding random noise) and the representations based on the final FPCA solution

 LV_{ψ} Metric $\in [0,1]$ quantifying the overlap of the simulated and

estimated FPC bases

 MISE_h Comparison of the simulated and estimated warping functions

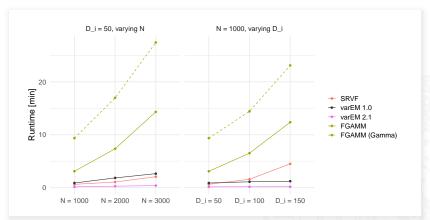
 MSE_d Comparison of the simulated and estimated domain lengths

back to sim study results

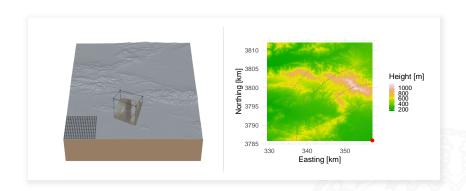
Simulation Study

Median runtimes

for one setting of the simulation study with amplitude rank 2-3 and no incompleteness, based on 20 runs for each parameter combination.

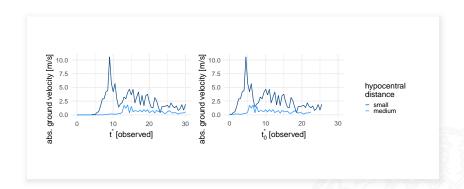


▶ back to sim study results



- Focus on wave propagation in northwest direction and close to the hypocenter
- Focus on t_0^* as the time since the arrival of seismic P-waves
- ⇒ Joint approach with Gamma distribution and trailing incompleteness

Seismic Application



- · Focus on wave propagation in northwest direction and close to the hypocenter
- Focus on t_0^* as the time since the arrival of seismic P-waves

⇒ Joint approach with Gamma distribution and trailing incompleteness

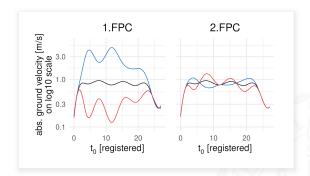
Seismic Application

Seismic application – Estimation details

- Used penalization parameter $\lambda = 0.004$
- 10 joint iterations, taking overall 3:31h, using a parallelized call for the registration steps with 5 cores
- The FPCs were chosen to explain 95% of amplitude variation

Curves and FPCs

with the first two FPCs visualized by the mean curve $+/-2\cdot\sqrt{\hat{ au}_k}\cdot\psi_k(t)$

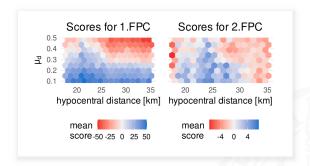


- \Rightarrow FPC 1 $\hat{=}$ overall magnitude with two peaks caused by surface waves
- \Rightarrow FPC 2 $\hat{=}$ salience of the initial peak

Seismic Application

Heatmaps of estimated phase and amplitude variation

conditional on hypocentral distance and the dynamic coefficient of friction μ_d



- \Rightarrow Overall ground motion shows strong association with μ_d
- \Rightarrow Initial peaks are most pronounced at hypocentral distances \sim 25km