

simm.fda package

a short description

Niels Olsen

June 2016

Latest updates: 22-3-22. Last updates were in 2016, this document is not yet up-to-date with the changes since then.

The `simm.fda` package is based on the article <https://doi.org/10.1111/rssc.12276> [O20] (preprint available at <http://arxiv.org/abs/1606.03295>), which will be general reference for modelling & methods. This document is meant as reference for the software.

The `simm.fda` package is based on the `pavpop` package by L.L. Raket (link: <https://github.com/larslau/pavpop>), which largely does the same thing as `simm.fda` in a one-dimensional setting. Some of the functions are taken directly from this package, but the two packages are not compatible and have significant differences in their internal structures.

The statistical setting for `simm.fda` is the model given by

$$\mathbf{y}_i(t_k) = \boldsymbol{\theta}_i(v_i(t_k)) + \mathbf{x}_i(t_k) + \boldsymbol{\varepsilon}_{ik} \quad (1)$$

where v_i is the unknown warp function, $\boldsymbol{\theta}_i : [0, 1] \rightarrow \mathbb{R}^d$ the expected trajectory for subject i and \mathbf{x}_i and $\boldsymbol{\varepsilon}_i$ represents amplitude covariance and noise. $\boldsymbol{\theta}_i$ can either be common to all subjects or subject to a design; $\boldsymbol{\theta}_i(t) = A_i\boldsymbol{\beta}(t)$.

It is assumed that $t \in [0, 1]$, although the model easily generalises to any domain $[a, b]$. However the software only accepts the domain $[0, 1]$, and hence the user must scale her time points to be within the domain $[0, 1]$. This may be changed in the future.

$\boldsymbol{\theta}$ (or the $\boldsymbol{\beta}$'s) is modelled as a spline, whose basis is assumed to be the same for all coordinates (and $\boldsymbol{\beta}$ s).

Estimation Estimation in (1) is done using localized iterative linearisation. The inner iteration consists of maximizing the posterior likelihood, i.e. predicting the warp given the current parameters, and estimating spline coefficients. The computationally intensive outer iteration consists of maximizing variance parameters in the linearized model. For the inner iterations warp change is used as criterion for convergence. The outer iterations will continue until the likelihood no longer decreases or maximum number of iterations has been reached.

Warp We have $v_i(t_k) := v(t, \mathbf{w}_i)$, where $w \in \mathbb{R}^d$ is a latent variable describing the warp. We shall refer to $v(\cdot, \cdot)$ as the *warp function*, and w as the (latent) *warp values*. Warping functions are assumed to be defined on the interval $[0, 1]$.

Warping functions are defined via `make_warp_fct`, default is to use type = 'smooth' with a small number of knots.

Basis function `simm.fda` supplies B-spline bases, Fourier bases and natural cubic splines ('ns'), defined via `make_basis_fct`. The user can provide her own basis functions.

Data The y 's have to be provided as list of $m_i \times K$ matrices, corresponding to a list of t 's. y may contain missing values as long as they are missing in all coordinates.

The amplitude covariances should be provided as function `function(t, param)` with the possibility of multiplying by a scale parameter, as a common scale parameter for warp and amplitude is always assumed. The result is a $Km_i \times Km_i$ matrix. The following order is used: The first indices correspond to the values/entries of the first coordinate, followed by values/entries of the second coordinate etc. Be aware of overparametrization, especially when using a dynamical covariance structure.

The warp covariances should be provided as functions `function(t, param)`. Do not use a scale parameter (unless the amplitude covariance is non-scalable), as a scale parameter is already a part of the model. Also warp covariances must have an attribute `param` that contains default parameter values. This requirement will be relaxed in the future.

Warp covariance may be null. In that case, all aspects in model that arise due to warp variation are set to zero.

Object returned by ppMulti `ppMulti` will return a list containing

- w : Estimated warp values
- `amp_cov_par`: Estimated amplitude covariance parameters.

- **warp_cov_par**: Estimated warp covariance parameters.
- **like**: Negative log-likelihood as estimated by the linearization.
- **sigma** Estimated σ in the model.
- **c**: Estimated spline coefficients. A matrix of dimension $\text{length}(\text{design}) \cdot K \times \text{df}(\text{spline})$. The estimated coefficients can be found using $(I_K \otimes \text{design}_i) \cdot c$.

Control parameters for ppMulti The program has spawned a multitude of control parameters. Most important are **iter**, the number of iterations.

[NEEDS UPDATE] Optimization parameters:

- **randomCycle** Two-vector. At each iteration after iteration **randomCycle**[2], randomly choose **randomCycle**[1] parameters for optimization. Overrides **paramMax**.
- **optimRule** Better choice than **randomCycle**. A list of vectors, where each vector contains parameters to be used in optimization, recycling vectors if necessary, e.g. **optimRule** = **list(c(1), c(2,3,4), c(5,6,7))** will optimize over parameter 1 in 1st iteration, parameters 2,3,4 in 2nd etc.
- **use.nlm**: Use **nlm** as opposed to **optim** for optimization? First index for outer loop, second index for inner loop. **nlm** might be better (or at least differently) behaved.

Dynamical covariance structure

Referring to Proposition 1 of [O20], dynamical covariance structure is a general framework for constructing covariance matrices using few parameters. It is given by the covariance function K constructed below:

Let f be a temporal covariance function, let $0 = s_1 < \dots < s_\ell = 1$ be anchor points, let $A_1, \dots, A_\ell \in \mathbb{R}^{q \times q}$ be a set of positive definite matrices. For each $t \in [0, 1]$ define $B_t \in \mathbb{R}^{q \times q}$ as the unique positive definite matrix satisfying

$$B_t^\top B_t = \frac{s_{k+1} - t}{s_{k+1} - s_k} A_k + \frac{t - s_k}{s_{k+1} - s_k} A_{k+1} \quad \text{for } t \in [s_k, s_{k+1}]. \quad (2)$$

For all $s, t \in [0, 1]$, define $K(s, t) = f(s, t) B_s^\top B_t \in \mathbb{R}^{q \times q}$.

Dynamical covariance structure is implemented in **simm.fda** using **kovMat**.

Timefunction concept A `timefunction` is a temporal covariance function on the domain $[0,1]$. `simm.fda` provides Matern covariance (including OU processes) and Brownian motion and bridge.

The user can provide her own time function. It has to be a function on the form `myfct(s,t, parameters) = Cov(X(s), X(t))` that can be vectorized by `outer`.

Time functions are passed to `kovMat` or used by calling `outer(t, t, mycft, ...)` using the syntax of `outer`.

Bridge concept A *bridge covariance* is defined by

$$f(s, t) = a + b(\min(s, t) - st)$$

where a is the 'stationary' part, and b the 'bridge' part.

This extends to any polynomial p , defining a polynomial bridge covariance by:

$$f(s, t) = p(\min(s, t) \cdot (1 - \max(s, t)))$$

Note that the bridge covariance is not the same construction as when conditioning a stochastic process X on its endpoint value.

Bridge covariances can (and should) be combined with other covariance functions by taking element-wise products, $A \circ B$. The validity of this construction follows from Schur's lemma.

Parallelization Parts of the estimation process can be parallelized, see the individual functions for details. Parallelization is done via `foreach`, and the user has to provide the parallel back-end herself. This is a deliberate design strategy.

Other parts of the software can be parallelized too, but these parts are not computationally heavy enough in order for me to believe parallelization has any advantage.

Data Example

The package includes a small data example. Data consists of four sequences of foot trajectories in x/y/z-space. The sequences are four repetitions of the same individual walking at a normal pace. It is part of much larger data set available from <http://mocap.cs.cmu.edu/>.

The variable `MCD.data` contains observation values, and `MCD.time` contains time points scaled to be within $[0.1, 0.9]$. A fairly simple model is fitted in the help page of `ppMulti`. Since fitting the model will generally take a long time, the output is available from `MCD.res`.

Functions

There are three 'main functions': `ppMulti`, `ppMulti.em` and `simfd.ed`.