



Function-on-Function Mixture of Experts Regression Models

Speaker: Jean Steve TAMO TCHOMGUI^{1,2}

work supervised by: Vincent BARRIAC¹, Guillaume FRAYSSE¹, Stéphane CHRETIEN², Julien JACQUES²

December 18th, 2022

¹ Orange Innovation, Paris France. ² Univ Lyon, Univ Lyon 2, ERIC, Lyon France.

Outline

- 1 Motivation
- 2 Function-on-Function Mixture Models
- 3 Inference
- 4 Simulation study
- 5 Future works

Motivations

- My PhD is entitled "Predicting the Quality of Experience of 5G verticals" from networks metrics;
- The shape of the available data led us to the Functional Data Analysis Framework;
- It further led us to the particular case of function-on-function regression models;
- With heterogeneity in the data, the latent class model is well adapted;
- For predictive modeling, the Mixture of Experts (MoE) is a more convenient architecture.

Data

Data ($1 \leq i \leq n$):

- Inputs: $X_i(t) = (X_i^1(t), \dots, X_i^p(t))$

- Output: $Y_i(t)$

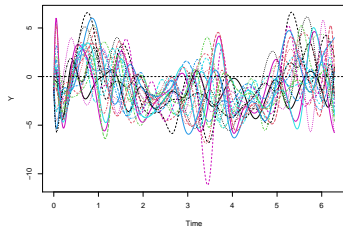
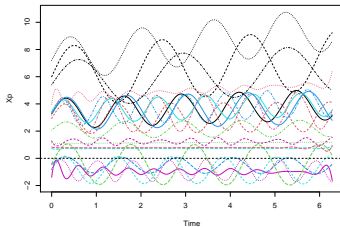
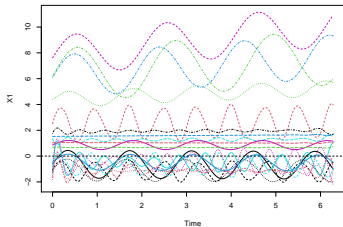


Table of Contents

- 1 Motivation
- 2 Function-on-Function Mixture Models**
- 3 Inference
- 4 Simulation study
- 5 Future works

Models for homogeneous data set

We consider the following models, Ramsay and Silverman (2005):

- **Concurrent model:**

$$Y_i(t) = \beta_0(t) + \sum_{l=1}^p \beta_l(t) X_i^l(t) + \varepsilon_i(t) = \beta(t)^\top X_i(t) + \varepsilon_i(t), \quad (1)$$

Models for homogeneous data set

We consider the following models, Ramsay and Silverman (2005):

- **Concurrent model:**

$$Y_i(t) = \beta_0(t) + \sum_{l=1}^p \beta_l(t) X_i^l(t) + \varepsilon_i(t) = \beta(t)^\top X_i(t) + \varepsilon_i(t), \quad (1)$$

- **(Feed forward) integral model:**

$$Y_i(t) = \gamma_0(t) + \sum_{l=1}^p \int_0^t \gamma_l(s, t) X_i^l(s) ds + \varepsilon_i(t) = \gamma_0(t) + \int_0^t \gamma(s, t)^\top X_i(s) ds + \varepsilon_i(t). \quad (2)$$

Mixture of Experts (MoE), Jacobs et al. (1991)

- **MoE** is Flexible supervised learning architecture dedicated for complex relationship;

Mixture of Experts (MoE), Jacobs et al. (1991)

- **MoE** is Flexible supervised learning architecture dedicated for complex relationship;
- The MoE regression model with functional response involving functional covariates in a K-component mixture is defined by:

$$f\left(Y_i(t) \mid X_i(t), t \in \mathbb{T}, \Psi\right) = \sum_{k=1}^K \underbrace{\pi_k(X_i(t), t \in \mathbb{T}, \alpha_k)}_{\text{component mixture weight}} \underbrace{f_k(Y_i(t) \mid X_i(t), \theta_k)}_{\text{Local expert } k} \quad (3)$$

Component mixture weights

- Model for the gated network function based on a functional multinomial logistic model. Dayton and Macready (1988), Chamroukhi et al. (2022):

$$\begin{aligned}\pi_k(X_i(t), t \in \mathbb{T}, \alpha_k) &= \mathbb{P}(Z = k | X_i(t), t \in \mathbb{T}, \alpha_k) \\ &= \frac{\exp(h_k(X_i(t), t \in \mathbb{T}, \alpha_k))}{1 + \sum_{k'=1}^{K-1} \exp(h_{k'}(X_i(t), t \in \mathbb{T}, \alpha_{k'}))}\end{aligned}\tag{4}$$

$$\text{Where } h_k(X_i(t), t \in \mathbb{T}, \alpha_k) = \alpha_{k,0} + \int_{\mathbb{T}} \alpha_k(t) X_i(t) dt$$

Locals experts

We use for local models the functional linear model **(1)** defined in an expert k by:

$$Y_i(t) = \beta_{0,k}(t) + \sum_{l=1}^p \beta_{l,k}(t) X_i^l(t) + \varepsilon_i(t) = \beta_k(t)^\top X_i(t) + \varepsilon_i(t) \quad \text{if } i \in \text{Group}_k$$

with $\beta_k(t) = (\beta_{0,k}(t) \quad \beta_{1,k}(t) \quad \dots \quad \beta_{p,k}(t))^\top$, $X_i(t) = (1 \quad X_i^1(t) \quad \dots \quad X_i^p(t))^\top$.

We expand in functional basis the functional covariates and parameters:

$$x_i^l(t) = \sum_{j=1}^{L_{x^l}} x_{ij}^l B_j^l(t) = B^l(t)^\top x_i^l \quad \text{and} \quad \beta_{l,k}(t) = \sum_{j=1}^{L_{\beta^l}} b_{j,k}^l \phi_j^l(t) = \phi^l(t)^\top b_k^l$$

The model can be written with this expression as linear mixed model:

$$Y_i(t) = b_k^\top \underbrace{\Phi(t)^\top B(t)}_{R_i(t)} x_i + \varepsilon_i(t) \tag{5}$$

Table of Contents

- 1 Motivation
- 2 Function-on-Function Mixture Models
- 3 Inference**
- 4 Simulation study
- 5 Future works

Estimation via the EM algorithm, Dempster et al. (1977)

We combine models (4) and (5) in K components mixture as:

$$f\left(Y_i(t) \mid X_i(t), t \in T, \Psi\right) = \sum_{k=1}^K \frac{\exp\left(\alpha_{k,0} + \delta_k^\top r_i\right)}{1 + \sum_{k'=1}^{K-1} \exp\left(\alpha_{k',0} + \delta_{k'}^\top r_i\right)} \Phi\left(Y_i(t); b_k^\top R_i(t), \sigma_k^2\right); \quad (6)$$

with $\Psi = \left((\alpha_{1,0}, \delta_1, b_1, \sigma_1^2), \dots, (\alpha_{K,0}, \delta_K, b_K, \sigma_K^2) \right)$ the parameters vector to estimate.

EM algorithm

- **E-step:** Given the current parameter Ψ^{old} , compute here the posterior probability that the individual i belongs to the class k :

$$\hat{p}_{ik} = \frac{\pi_k(X_i(t), t \in T, \alpha_k^{\text{old}}) \Phi(Y_i(t); b_k^{\text{old}\top} R_i(t), \sigma_k^{2\text{old}})}{\sum_{k'=1}^K \pi_{k'}(X_i(t), t \in T, \alpha_{k'}^{\text{old}}) \Phi(Y_i(t); b_{k'}^{\text{old}\top} R_i(t), \sigma_{k'}^{\text{old}^2})}$$

EM algorithm

- **E-step:** Given the current parameter Ψ^{old} , compute here the posterior probability that the individual i belongs to the class k :

$$\hat{p}_{ik} = \frac{\pi_k(X_i(t), t \in \mathbb{T}, \alpha_k^{\text{old}}) \Phi(Y_i(t); b_k^{\text{old}\top} R_i(t), \sigma_k^{2\text{old}})}{\sum_{k'=1}^K \pi_{k'}(X_i(t), t \in \mathbb{T}, \alpha_{k'}^{\text{old}}) \Phi(Y_i(t); b_{k'}^{\text{old}\top} R_i(t), \sigma_{k'}^{2\text{old}})}$$

- **M-step:** Given the posterior probabilities, we update the parameters by maximizing:

$$Q(\Psi^{\text{new}} | \Psi^{\text{old}}) = Q_1(b^{\text{new}}, \sigma^{2\text{new}} | \Psi^{\text{old}}) + Q_2(\delta^{\text{new}} | \Psi^{\text{old}})$$

Where:

$$Q_1(b^{\text{new}}, \sigma^{2\text{new}} | \Psi^{\text{old}}) = \sum_{i=1}^N \sum_{k=1}^K \hat{p}_{ik} \log \left(\Phi(Y_i(t); b_k^{\text{new}\top} R_i(t), \sigma_k^{2\text{new}}) \right)$$

$$Q_2(\delta^{\text{new}} | \Psi^{\text{old}}) = \sum_{i=1}^N \sum_{k=1}^K \hat{p}_{ik} \log \left(\pi_k(X_i(t), t \in \mathbb{T}, \alpha_k^{\text{new}}) \right)$$

EM algorithm

- Q_1 and Q_2 can be maximized separately using weighted ML estimation of GLMs and multinomial logistic model respectively;
- The algorithm is stopped if the relative change in the log-likelihood is smaller than a pre-specified ϵ ;
- We repeat the EM algorithm with different initializations and choose the one with the maximum likelihood;
- Implementations details can be seen in Grün and Leisch (2008).

Model selection:

- The number K of mixture components can be chosen using BIC.

Table of Contents

- 1 Motivation
- 2 Function-on-Function Mixture Models
- 3 Inference
- 4 Simulation study**
- 5 Future works

Simulation study

- We simulate two data sets with $n = 500$ and $n = 1000$ each with $p = 5$ functional predictors defined on $[0, 2\pi]$ given by the equation:

$$\begin{cases} U_i^l(t_j) = \xi_{i,1}^l + \left(\log(10 + t_j)\right)^{\xi_{i,2}^l} + \xi_{i,3}^l \sin\left(\frac{2\pi t_j}{\xi_{i,4}^l}\right) \\ \xi_{i,r}^l \quad \text{constants drawn from } \mathcal{U}([-1, 1]), \quad 1 \leq r \leq 4 \end{cases}$$

- And the functional response by the concurrent model with the formula:

$$Y_i(t_j) = \beta_{0,k}(t_j) + \sum_{l=1}^p \beta_{l,k}(t_j) X_i^l(t_j) + \varepsilon_{ij} \quad \text{if } i \in \text{Group}_k,$$

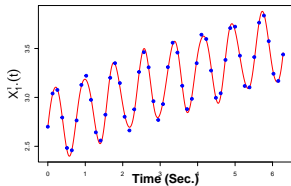
$$\text{with } \begin{cases} \beta_{0,k}(t_j) = \left(\log(10 + t_j)\right)^{\rho_{0,k}} \\ \beta_{l,k}(t_j) = \rho_1^{l,k} \sin\left(\frac{2\pi t_j}{\rho_2^{l,k}}\right) \end{cases} \quad \text{where } \begin{cases} \varepsilon_{ij} \sim \mathcal{N}(0, 4), \\ \rho_{0,k}, \rho_1^{l,k}, \rho_2^{l,k} \text{ are chosen constants.} \end{cases}$$

- $K = 4$ groups, so we have $4 \times 6 = 24$ parameters ; In each scenario, we run $N = 50$ experiences.

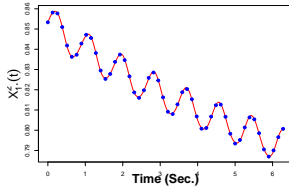
Functional predictors and functional concurrent response for a randomly chosen individual:

The blue dots are the observed discrete data and the red curves is the functional reconstructions of predictors

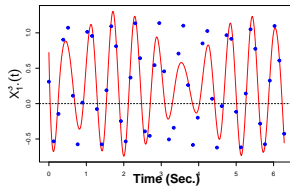
Predictor 1



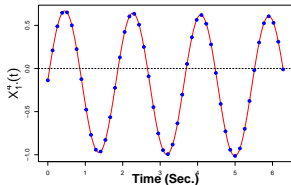
Predictor 2



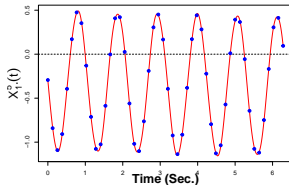
Predictor 3



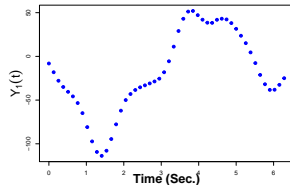
Predictor 4



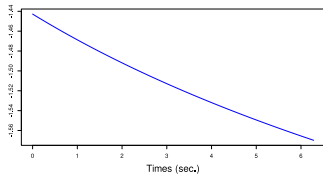
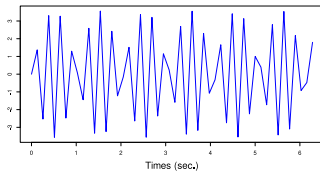
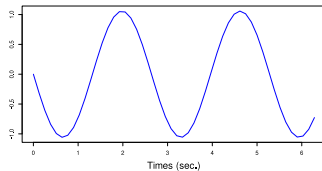
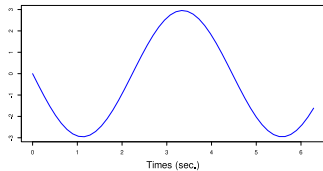
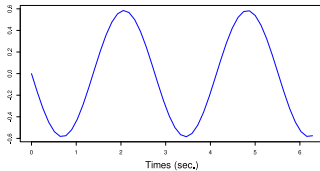
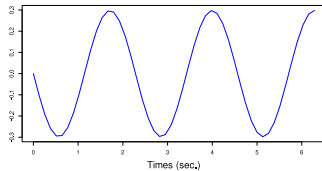
Predictor 5



Response



Functional parameters for the component (expert) 1:

 $\beta_{0,1}(t)$  $\beta_{1,1}(t)$  $\beta_{2,1}(t)$  $\beta_{3,1}(t)$  $\beta_{4,1}(t)$  $\beta_{5,1}(t)$

Assessment of quality of estimation

- **Choice of K :** Ratio of models with the right number of components ;
- **Estimation of the $\beta_l(t)$:** The accuracy of the estimated parameters of models that have the good number of components via the Root Mean Square Error (RMSE):

$$\text{RMSE}(\beta_l(\cdot)) = \left[\frac{1}{m} \sum_{j=1}^m (\beta_l(t_j) - \hat{\beta}_l(t_j))^2 \right]^{1/2}.$$

- **Prediction quality:** via the Mean Relative Prediction Error (MRPE) given by:

$$\text{MRPE} = \frac{1}{m} \sum_{j=1}^m \left(\frac{\sum_{i=1}^n (Y_i(t_j) - \hat{Y}_i(t_j))^2}{\sum_{i=1}^n Y_i(t_j)^2} \right).$$

Choice of K

We evaluate the efficiency of BIC for choosing K :

- For each sample, different values for $K \in \{2, \dots, 6\}$ has been tested, and BIC is used to select the best one;
- **Ratio of models with the right number of components:**

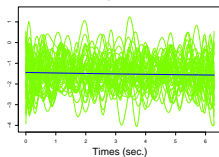
Number of classes		3	4	5
% of models	n.train = 500	6%	86%	8%
	n.train = 1000	0%	92%	8%

Table: Results of our $N = 50$ simulations in terms of number of components.

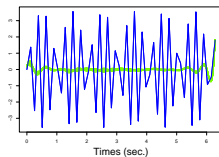
Estimation of the of the $\beta_l(t)$

Figure: Estimated (green) and actual (blue) parameters for our simulations with $n_{train} = 500$.

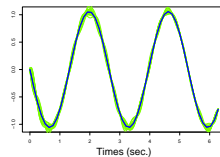
Component 1



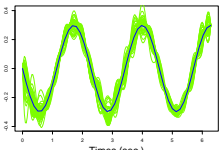
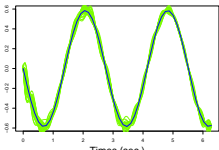
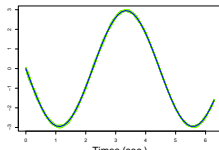
$\beta_{0,1}(t)$



$\beta_{1,1}(t)$

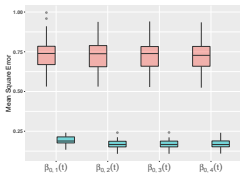


$\beta_{2,1}(t)$

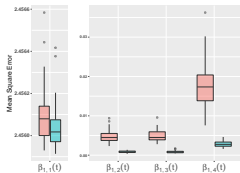


Estimation of the $\beta_l(t)$

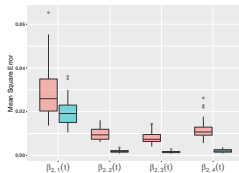
Figure: Boxplot of RMSE in all the estimated parameters for our two scenarios : $n = 500$ (red) and $n = 1000$ (blue).



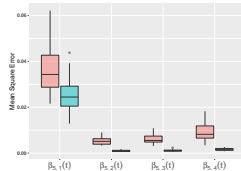
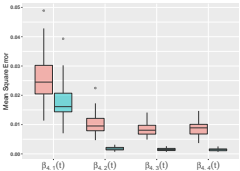
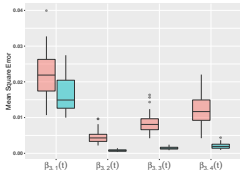
$\beta_0(t)$



$\beta_1(t)$

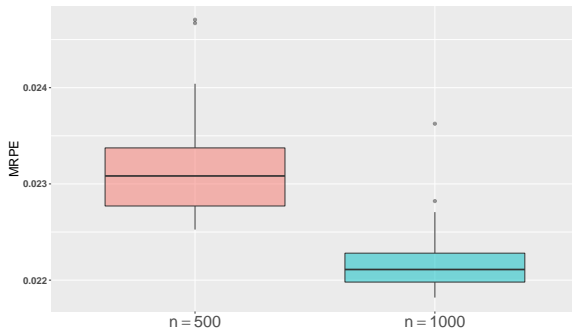


$\beta_2(t)$

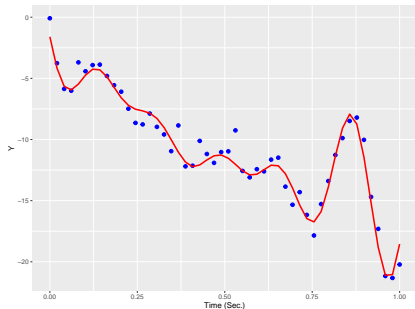


Prediction quality : MRPE

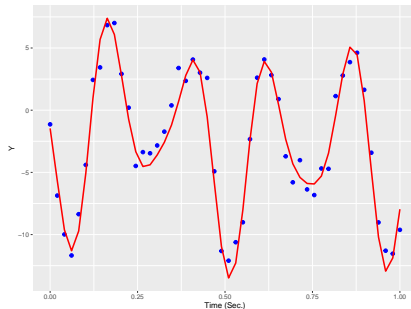
We evaluate here the prediction error on a test sample of length $n_{test} = 2000$ through the MRPE :



Prediction Quality : some plots



(a) prediction 1



(b) prediction 2

Figure: Prediction and actual values of the functional response for two randomly chosen individuals.

Conclusion and Future works

- MoE is a flexible architecture that allows to efficiently represents complex relationship ;
- EM algorithm designed here for functional data can be seen as extension vectorized version of MoE ;

Conclusion and Future works

- MoE is a flexible architecture that allows to efficiently represents complex relationship ;
- EM algorithm designed here for functional data can be seen as extension vectorized version of MoE ;
- With the functional basis expansion of our parameters, we must enhance both the local expert and gated network estimated parameters by applied a roughness penalty to get them interpretable :

$$\text{Pen}(\beta_l) = \lambda_l \int \beta_l''(t)^2 dt = \lambda_l \int \left[\sum_{j=1}^{L_{\beta^l}} b_j^l \phi_j^l(t) \right]^2 dt = \lambda_l \sum_{s,k=1}^{L_{\beta^l}} b_s^l b_k^l \Phi_{sk}^l$$

$$\text{with } \Phi_{sk}^l = \int \phi_s^l(t) \phi_k^l(t) dt.$$

Some references

- F. Chamroukhi, N. T. Pham, V. H. Hoang, and G. J. McLachlan, “Functional mixtures-of-experts,” arXiv preprint arXiv:2202.02249, 2022.
- C. M. Dayton and G. B. Macready, “Concomitant-variable latent-class models,” Journal of the american statistical association, vol. 83, no. 401, pp. 173–178, 1988.
- A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” Journal of the Royal Statistical Society: Series B (Methodological), vol. 39, no. 1, pp. 1–22, 1977.
- B. Grün and F. Leisch, “Flexmix version 2: Finite mixtures with concomitant variables and varying and constant parameters,” Journal of Statistical Software, vol. 28, no. 4, p. 1–35, 2008. [Online]. Available: <https://www.jstatsoft.org/index.php/jss/article/view/v028i04>
- R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, “Adaptive Mixtures of Local Experts,” Neural Computation, vol. 3, no. 1, pp. 79–87, 03 1991. [Online]. Available: <https://doi.org/10.1162/neco.1991.3.1.79>
- J. Ramsay and B. W. Silverman, Functional Data Analysis, 2nd ed., ser. Springer Series in Statistics. New York: Springer-Verlag, 2005. [Online]. Available: <https://www.springer.com/gp/book/9780387400808>

Questions?

Appendix - Recover the functional nature of the predictors

- In the particular case of the cubic spline model, the expression in functional basis $\{B_l(t)\}_{l \geq 1}$ truncated at q becomes $\{1, t, t^2, t^3, (t - \tau_1)_+^3, \dots, (t - \tau_{q-4})_+^3\}$ and then

$$x_i^l(t) = x_{i1}^l + x_{i2}^l t + x_{i3}^l t^2 + x_{i4}^l t^3 + \sum_{l=1}^{q-4} x_{i(4+l)}^l (t - \tau_l)_+^3.$$

$$\begin{aligned} \begin{pmatrix} x_i(t_{i1}) \\ x_i(t_{i2}) \\ \vdots \\ x_i(t_{im_i}) \end{pmatrix} &= \begin{pmatrix} x_{i1} + x_{i2} t_{i1} + x_{i3} t_{i1}^2 + x_{i4} t_{i1}^3 + \sum_{l=1}^{q-4} x_{i(4+l)} (t_{i1} - \tau_l)_+^3 \\ x_{i1} + x_{i2} t_{i2} + x_{i3} t_{i2}^2 + x_{i4} t_{i2}^3 + \sum_{l=1}^{q-4} x_{i(4+l)} (t_{i2} - \tau_l)_+^3 \\ \vdots \\ x_{i1} + x_{i2} t_{im_i} + x_{i3} t_{im_i}^2 + x_{i4} t_{im_i}^3 + \sum_{l=1}^{q-4} x_{i(4+l)} (t_{im_i} - \tau_l)_+^3 \end{pmatrix} \\ &= \begin{pmatrix} 1 & t_{i1} & t_{i1}^2 & t_{i1}^3 & (t_{i1} - \tau_1)_+^3 & (t_{i1} - \tau_2)_+^3 & \dots & (t_{i1} - \tau_{q-4})_+^3 \\ 1 & t_{i2} & t_{i2}^2 & t_{i2}^3 & (t_{i2} - \tau_1)_+^3 & (t_{i2} - \tau_2)_+^3 & \dots & (t_{i2} - \tau_{q-4})_+^3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_{im_i} & t_{im_i}^2 & t_{im_i}^3 & (t_{im_i} - \tau_1)_+^3 & (t_{im_i} - \tau_2)_+^3 & \dots & (t_{im_i} - \tau_{q-4})_+^3 \end{pmatrix} \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iq} \end{pmatrix} \end{aligned}$$

$$x_i = \mathcal{B}_i x_i.$$

It is recommended (Ruppert, 2002) to choose q large enough to capture the complexity of the predictor and then apply a roughness penalty to limit overfitting of the form:

$$\text{(C1)} \max_{5 \leq l \leq q} |x_{il}| < C \quad \text{(C2)} \sum_{l=5}^q |x_{il}| < C \quad \text{(C3)} \sum_{l=5}^q x_{il}^2 < C.$$

$$\begin{cases} \min_{x_i} \|x_i - \mathcal{B}_i x_i\|^2 \\ \text{s.c. } x_i^\top D x_i \leq C \end{cases} \iff \min_{x_i} \|x_i - \mathcal{B}_i x_i\|^2 + \lambda x_i^\top D x_i \quad (\lambda \geq 0)$$

$$\text{with } D = \begin{bmatrix} 0_{4 \times 4} & 0_{4 \times (q-4)} \\ 0_{(q-4) \times 4} & \mathbf{I}_{(q-4) \times (q-4)} \end{bmatrix}$$

Appendix - Penalizing the model

- The estimation of functional parameters $\beta(t)$ expressed in functional basis depends on the number of basis functions L_β ;
- Unfortunately, the functional basis based fitting procedure lead to a non regular estimate curve $\hat{\beta}(t)$ and then hard to interpret;
- To correct this non-regularity effect, we enhance the model by introducing a Ridge roughness penalty in the model.

In this case, the penalized least squares risk to be minimized is given by:

$$\mathcal{R}(b) = \sum_{i=1}^n \sum_{j=1}^m \left(Y_i(t_j) - b^\top R_i(t_j) \right)^2 + \sum_{l=0}^p \text{Pen}(\beta_l), \quad (7)$$

where $\text{Pen}(\beta_l) = \lambda_l \int \beta_l''(t)^2 dt = \lambda_l \int \left[\sum_{j=1}^{L_{\beta^l}} b_j^l \phi_j^l(t) \right]^2 dt = \lambda_l \sum_{s,k=1}^{L_{\beta^l}} b_s^l b_k^l \Phi_{sk}^l$ with $\Phi_{sk}^l = \int \phi_s^l(t) \phi_k^l(t) dt$.

Appendix - Penalizing the model

Return to presentation

Then for a fixed value λ_l , the estimation of $(\beta_l(t))_{0 \leq l \leq p}$ is calculated by solving the problem:

$$\begin{aligned} \min_b \mathcal{R}(b) &= \min_b \sum_{i=1}^n \sum_{j=1}^m \left(Y_i(t_j) - b^\top R_i(t_j) \right)^2 + \sum_{l=0}^p \lambda_l \sum_{s,k=1}^{L_{\beta^l}} b_s^l b_k^l \Phi_{sk}^l \\ &= \min_b (Y - Rb)^\top (Y - Rb) + b^\top (\lambda P) b, \end{aligned}$$

where λP the matrix of dimension $L_{\beta} \times L_{\beta}$ given by:

$$\lambda P = \begin{pmatrix} \lambda_0 \Psi^0 & \mathbb{0} & \dots & \mathbb{0} \\ \mathbb{0} & \lambda_1 \Psi^1 & \dots & \mathbb{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{0} & \mathbb{0} & \dots & \lambda_p \Psi^p \end{pmatrix} \quad \text{with } \Psi^l = \begin{pmatrix} \Phi_{11}^l & \Phi_{12}^l & \dots & \Phi_{1L_{\beta^l}}^l \\ \Phi_{21}^l & \Phi_{22}^l & \dots & \Phi_{2L_{\beta^l}}^l \\ \vdots & \vdots & \ddots & \vdots \\ \Phi_{L_{\beta^l}1}^l & \Phi_{L_{\beta^l}2}^l & \dots & \Phi_{L_{\beta^l}L_{\beta^l}}^l \end{pmatrix}.$$

With $\mathbb{0}$ a notation to refer to the corresponded block of null matrix. As Ψ^l for any $0 \leq l \leq p$ is a symmetric positive-definite matrix, we can easily find its Cholesky decomposition which will be used for the practical implementation.