

¹ **Functional error modeling for uncertainty**
² **quantification in hydrogeology**

L. Josset,¹ D. Ginsbourger,² I. Lunati,¹

Corresponding author: L. Josset, ISTE, University of Lausanne, Switzerland. (laureline.josset@unil.ch)

¹ISTE, University of Lausanne,
Switzerland

²IMSV, University of Bern, Switzerland

Abstract.

Approximate models (proxies) can be employed to reduce the computational costs of estimating uncertainty. The price to pay is that the approximations introduced by the proxy model can lead to a biased estimation. To avoid this problem and ensure a reliable uncertainty quantification, we propose to combine Functional Data Analysis and Machine Learning to build error models that allow us to obtain an accurate prediction of the exact response without solving the exact model for all realizations. We build the relationship between proxy and exact model on a learning set of geostatistical realizations for which both exact and approximate solvers are run. Functional principal components analysis (FPCA) is used to investigate the variability in the two sets of curves and reduce the dimensionality of the problem while maximizing the retained information. Once obtained, the error model can be used to predict the exact response of any realization on the basis of the sole proxy response. This methodology is purpose-oriented as the error model is constructed directly for the quantity of interest, rather than for the state of the system. Also, the dimensionality reduction performed by FPCA allows a diagnostic of the quality of the error model to assess the informativeness of the learning set and the fidelity of the proxy to the exact model. The possibility of obtaining a prediction of the exact response for any newly generated realization suggests that the methodology can be effectively used beyond the context of uncertainty quantification, in particular for Bayesian inference and optimization.

1. Introduction

26 The major challenge in hydrogeology is to deal with an incomplete knowledge of aquifer
27 properties, which are usually measured only at few, discrete locations. This lack of infor-
28 mation makes it impossible to address hydrogeological problems in a deterministic sense.
29 The problem is typically stated in a stochastic framework and Monte Carlo simulations
30 are used to propagate the uncertainty on aquifer properties to the quantities of interest
31 [*Dagan, 2002*]. A typical example is the prediction of the fate of a contaminant, which de-
32 pends on the heterogeneity structure of the aquifer. The uncertainty on the contaminant
33 breakthrough curve at a given location is estimated by solving the transport problem in
34 a set of realizations, which represent the uncertainty on the permeability of the aquifer.
35 The ensemble of the responses in the different realizations provides a sample of reference
36 of the breakthrough curves.

37 Despite the appealing conceptual simplicity of this approach, problems arise when many
38 realizations have to be considered and a large number of expensive flow and transport sim-
39 ulations have to be performed: computational cost quickly becomes prohibitive. To avoid
40 this computational bottleneck, the problem is approximated either by coarsening the
41 description of aquifer properties (standard upscaling techniques can be used to this end
42 [*Wen and Gómez-Hernández, 1996; Renard and de Marsily, 1997; Christie, 1996; Durlof-*
43 *sky, 2005*]) or by simplifying the description of the physical processes, thus employing an
44 approximate model or proxy (e.f., *Scheidt and Caers [2009a]*).

45 The price to pay for these simplifications is that inference based on the computed
46 responses could lead to a wrong uncertainty quantification. If the approximation is phys-

47 ically motivated, the bias can be safely ignored. Effective computational gains, however,
48 usually require very crude approximations whose effects on the uncertainty quantification
49 is difficult to assess beforehand. To avoid this problem, the proxies are typically employed
50 only to identify a representative subset of realizations for which the exact model is solved.
51 This is the strategy of ranking methods [*McLennan and Deutsch, 2005; Ballin et al.,*
52 *1992*], or distance kernel methods [*Scheidt and Caers, 2009a*]. In such case, it is crucial
53 to evaluate to which extent the proxy is informative of the exact model response.

54 While it is generally acknowledged that an error analysis is necessary [*Christie et al.,*
55 *2005*], it is rarely performed. Although approaches that entail a systematic analysis and
56 the construction of error models have been applied to flow in porous media (e.g., to cor-
57 rect fluid-properties approximations *O'Sullivan and Christie [2005, 2006]* or approximate
58 numerical solvers *Josset and Lunati [2013]*), in most cases the appraisal of approximate
59 methods is performed for a very limited number of test cases, and it is assumed that they
60 behave similarly for a wider range of applications. This approach is not exempt from
61 problems because the informativeness of the proxy also depends on flow regimes and on
62 the specific quantities of interest.

63 In this paper, we propose a novel methodology to systematically build statistical error
64 models that describe the discrepancy between exact and approximate responses. Once
65 the error model is constructed, it is used to correct the approximate responses and predict
66 the responses expected from the exact model for all realizations. A characteristic of our
67 approach is that the error model is purpose oriented, that is, it is established directly for
68 the quantities of interest (in our case the breakthrough curve of a contaminant) and not
69 for the state of the system (for instance, the full saturation -or concentration- and pressure

70 fields). This reduces the complexity of the data to be handled (e.g., time-dependent curves
71 rather than time-dependent fields) while retaining all the relevant information.

72 Despite some similarities with the error models proposed by *Josset and Lunati* [2013],
73 two additional key features characterize the present approach: the description of sparse
74 data as continuous variables (time-dependent breakthrough curves), and the reduced di-
75 mensionality of the problem that is solved to construct the error model. To this end
76 we employ Functional Principal Component Analysis (FPCA [*Henderson, 2006*]), which
77 is a functional extension of PCA. The theoretical background is provided by Functional
78 Data Analysis (FDA), a discipline that gathers mathematical tools to construct and treat
79 continuous data. The description of continuous variables from sparse data is a problem
80 faced in many fields of research and not only in environmental applications. While func-
81 tional analysis is well established, FDA has been integrated as a whole only recently and
82 promoted by *Ramsay* [2006]; *Ramsay et al.* [2009]. It has since been applied in various
83 areas such as biomedical science, biomechanics, medicine or linguistic among others. We
84 refer to *Ullah et al.* [2013] for a recent review of the application of FDA over the last 20
85 years. More specifically to the domain of groundwater protection problem, FPCA has
86 been applied to interpret various contaminant concentrations in river quality [*Henderson,*
87 2006].

88 The paper is organized as follows. After a general problem statement (Sec. 2), we
89 introduce the formalism used and describe the methodology in detail (Sec. 3). Then, the
90 methodology is evaluated for a synthetic test case that represents a typical groundwater
91 problem (Sec. 4). The paper ends with a discussion of the performance and of prospective
92 applications (Sec. 5).

2. Problem statement

93 We consider a contamination problem in which a non-aqueous phase liquid (NAPL) is
94 accidentally released and forms a plume that contaminates the fresh water. We are inter-
95 ested in predicting the breakthrough curve of the pollutant at a given location (typically a
96 drinking well or a river that can be contaminated). Examples of NAPL contamination are
97 hydrocarbons spills, or leakage of chlorinated solvents such as TCE. As the NAPL is not
98 miscible with water and forms a separate phase, the evolution of the contamination plume
99 is governed by a set of nonlinear transport equations (Appendix A), which complicates
100 both the contaminant behaviour and the numerical resolution of the equations.

101 Due to sparse measurements, the properties of the aquifer are only partially known.
102 Their uncertainty is represented by a set of N_r geostatistical realizations of the per-
103 meability and porosity fields $\{R_i\}_{i=1,\dots,N_r}$. In brute force Monte Carlo approaches, this
104 uncertainty is propagated by solving the nonlinear multiphase transport model (here-
105 after “exact model”) and computing the NAPL breakthrough curve in each realization.
106 Here it is assumed that the resulting set of curves, $\{y_i(t)\}_{i=1,\dots,N_r}$, provides an accurate
107 representation of the uncertainty on the travel time.

108 Our goal is to find an approximation of the uncertainty without computing the full
109 set of exact curves $\{y_i(t)\}_{i=1,\dots,N_r}$. To this end we use a simplified model based on
110 the linear single-phase transport equations (hereafter “approximate model” or “proxy”),
111 which allows a relatively inexpensive calculation of the approximate breakthrough curves,
112 $\{x_i(t)\}_{i=1,\dots,N_r}$. To provide an accurate approximation of the uncertainty, we need to learn
113 the relationship between the proxy and the exact responses, such that an exact response
114 can be predicted from each proxy response.

115 We formulate this step in a standard machine learning framework: a statistical model
 116 relating the exact response curves (treated as outputs of the statistical model) to the
 117 proxy response curves (treated as inputs of the statistical model) is postulated. The
 118 parameters are estimated based on a learning set (or training set), i.e., a collection of pairs
 119 of response curves obtained with the two models for $N_l < N_r$ geostatistical realizations,
 120 $\{(x_i(t), y_i(t))\}_{i=1, \dots, N_l}$.

121 The statistical model relating the two sets of response curves (exact and proxy) is here
 122 restricted to the class of functional linear models [Ramsay, 2006], in which the relation-
 123 ships between the responses is

$$y_i = T(x_i) + \varepsilon_i \quad i \in [1, \dots, N_r], \quad (1)$$

124 where T is a bounded linear operator from the Hilbert space L_2 to itself, and the error
 125 functions ε_i are centered, independent, and typically assumed to meet further technical
 126 conditions [Cuevas *et al.*, 2002].

127 Since the identification of such statistical model is ill-posed, in practice further restric-
 128 tions on the form of T are made introduced to enable inferring T from the learning set.
 129 Two methods are suggested by Ramsay [2006]; Ramsay *et al.* [2009]: the full functional
 130 regression model and the Concurrent model. The full functional regression model allows
 131 capturing complex behaviours, but it is costly and requires the fine tuning of several
 132 smoothing parameters. The Concurrent model consists of a simpler functional linear
 133 regression. This method is fast, but quite rudimentary because the model uses only con-
 134 current features of the curves (additional details about the two models can be found in
 135 Appendix B).

136 In this paper, we follow a slightly different strategy: we appeal to a spectral approach
 137 and decompose the elements of the learning set on two *ad hoc* bases, one for the proxy
 138 and one for the exact responses. The response curves are then described in two spaces
 139 of dimensions $D_{ex} < N_l$ for the exact responses and $D_{app} < N_l$ for the proxy responses.
 140 A statistical model is constructed to relate the coefficients of the elements of one space,
 141 $y_i(t)$, to the coefficients of the elements of the other space, $x_i(t)$, as illustrated in Fig. 1.

Once the approximation \hat{T} of T is obtained from the learning-set, it is used to predict the exact responses of all realizations from of the approximate responses, i.e.,

$$\{\hat{y}_i = \hat{T}(x_i)\}_{i=1,\dots,N_r}, \quad (2)$$

142 and the uncertainty is quantified from the ensemble of predicted curves.

3. Methodology

143 The construction of the error model consists of four steps: first, functional objects
 144 are built from the data in the learning set; second, the dimensionality of the problem is
 145 reduced by decreasing the dimensions of the two functional spaces; third, the relationship
 146 between the approximate and exact responses is constructed; fourth, the error model is
 147 used to predict the exact responses from the proxy responses. These steps are illustrated
 148 in the flowchart in Fig. 2.

3.1. Recasting discretized curves as functional data

149 Both exact and proxy responses are obtained from numerical simulations and are rep-
 150 resented by contaminant breakthrough curves defined at discrete times. Therefore, we
 151 recast the time-discrete curves into time-continuous functions. This has two practical ad-
 152 vantages: first, it allows us to use the formalism of functional data analysis and the tools

153 that have been developed in this context; second, it permits to work with asynchronous
154 information about the curves, i.e., curves that have been sampled at different times. Note
155 that this step is essential in applications in which analytic solutions are used as proxies
156 or if the exact responses are provided by field measurements, which are typically acquired
157 with different temporal resolution.

158

Many functional bases are available to recast discretized curves into functional data. Here, we use a K -dimensional B-spline basis denoted by $\{\varphi_k(t)\}_{k \in [1, K]}$. To determine the coefficients, a linear combination of the elements of this basis is fitted to the data, which are represented as time dependent functions of the form

$$f(t) = \sum_{k=1}^K c_k \varphi_k(t) \quad (3)$$

159 *Ramsay* [2006] suggests two strategies to choose the basis and fit the coefficients to
160 data: either a low-dimension basis is used and the data are plainly projected (e.g., by
161 ordinary least squares), or a high-dimension basis is used with a roughness penalty to
162 mitigate overfitting. Both strategies allow not only to distinguish noise from information
163 but also to impose various constraints on the functional objects, e.g. positivity and/or
164 monotonicity. As our data (contaminant breakthrough curves) are typically fairly smooth,
165 a standard B -spline basis of small dimension can be used. We refer the readers to [*Ramsay*,
166 2006; *Ramsay et al.*, 2009] for more details about the notions of roughness penalty and
167 incorporation of constraints.

3.2. Functional reduction of the dimensionality

168 The previous step allows representing each exact response and each proxy response
 169 as a continuous function, i.e., $y_i(t)$ and $x_i(t)$, respectively. To decrease the dimension
 170 of the response spaces and the size of the regression problem, we employ Functional
 171 Principal Component Analysis, which is a functional extension of standard PCA and
 172 allows highlighting the main modes of variability in a sample of functions. Beside a small
 173 computational advantage, using spaces of lower dimension reduces the risk of over-fitting
 174 and allows us to visualize the data to assess the informativeness of the proxy response
 175 with respect to the exact response.

We apply FPCA to the exact and proxy responses in the learning set. Given the sample of proxy functions in the learning set, $\{x_i(t)\}_{i=1,\dots,N_l}$, with average $\bar{x}(t) = \frac{1}{N_l} \sum_{i=1}^{N_l} x_i(t)$ and estimated covariance function

$$\nu(t', t) = \frac{1}{N_l - 1} \sum_{i=1}^{N_l} [x_i(t') - \bar{x}_i(t')][x_i(t) - \bar{x}_i(t)], \quad (4)$$

FPCA constructs a non increasing sequence of eigenvalues of the estimated covariance function, $\mu_1^\circ \geq \mu_2^\circ \geq \dots \geq \mu_{N_l-1}^\circ$, by solving the functional eigenequation

$$\int \nu(t', t) \zeta_i^\circ(t) dt = \mu_i^\circ \zeta_i^\circ(t'). \quad (5)$$

The sequence of eigenfunctions (or harmonics) of the covariance function, $\{\zeta_1^\circ, \dots, \zeta_{N_l-1}^\circ\}$, satisfies the condition

$$\int \zeta_i^\circ(t) \zeta_j^\circ(t) dt = \delta_{ij}, \quad (6)$$

176 (where δ_{ij} is the Kronecker delta), and, together with the average $\bar{x}(t)$, form an or-
 177 thonormal basis for the space of the sampled approximate responses. The eigenvalue μ_i is
 178 also denoted as the *probe score variance* and the eigenfunction $\zeta_i^\circ(t)$ as *harmonic* [Ramsay

179 *et al.*, 2009]. The dimensionality of the response space can be optimally reduced consider-
 180 ing only the first D_{ex} and D_{app} for the exact response space and the proxy response space,
 181 respectively. The fact that the sequence of eigenvalues is non increasing guarantees that
 182 no other basis of size D_{app} can describe better the data; the total squared error introduced
 183 by discarding the eigenfuncions $(\zeta_i^\circ(t))_{i>D_{app}}$ is $\sum_{i=D_{app}+1}^{N_i-1} \mu_i^\circ$.

184 The basis allows us to approximate each proxy response as

$$x_i(t) \approx \tilde{x}_i(t) = \bar{x}(t) + \sum_{j=1}^{D_{app}} b_{ij}^\circ \zeta_j^\circ(t) \tag{7}$$

where

$$b_{ij}^\circ = \int [\bar{x}(t) - x_i(t)] \zeta_j^\circ(t) dt \tag{8}$$

185 is the projection of the deviation from the mean of the i^{th} approximate curve on the j^{th}
 186 harmonic ($\tilde{x}_i(t)$ denotes the approximation of $x_i(t)$ in terms of the first D_{app} harmonics).
 187 As in standard PCA, these coefficients are typically referred to as *scores*.

Although it offers an optimal dimensionality reduction with respect to the total mean squared error, the orthonormal basis might not be ideal to represent the information. The *varimax* algorithm [Kaiser, 1958] can be applied to find a suitable rotation that improve data interpretation while preserving the optimality of the result in terms of explained variance [Richman, 1986; Ramsay *et al.*, 2009]. Therefore, without any further loss of information, the approximate curves can be written as

$$\tilde{x}_i(t) = \bar{x}(t) + \sum_j^{D_{app}} b_{ij} \zeta_j(t), \tag{9}$$

where

$$b_{ij} = \int [\bar{x}(t) - x_i(t)] \zeta_j(t) dt \tag{10}$$

188 is the projection of the deviation from the mean of the i^{th} curves on the rotated harmonic
 189 $\zeta_j(t)$.

An analogous procedure is applied to the sample of exact responses in the learning set, $\{y_i(t)\}_{i=1,\dots,N_l}$, which is approximated as

$$\tilde{y}_i(t) = \bar{y}(t) + \sum_j^{D_{ex}} c_{ij} \eta_j(t), \quad (11)$$

where $\bar{y}(t)$ is the average, $\eta_j(t)$ the j^{th} harmonic of the (varimax) rotated orthonormal basis $\{\eta_i(t)\}_{i=1,\dots,D_{ex}}$, and

$$c_{ij} = \int [y_i(t) - \bar{y}(t)] \eta_j(t) dt \quad (12)$$

190 the score with respect to $\eta_j(t)$. (As for the proxy curve, the *tilde* denotes the restriction
 191 to the first D_{ex} harmonics).

3.3. Regression and error model

192 Once the problem dimensionality has been reduced by FPCA, we investigate the rela-
 193 tionships between the two sets of curves in the learning set approximated by considering
 194 the first D_{app} and D_{ex} harmonics, $\{\tilde{x}_i(t), \tilde{y}_i(t)\}_{i=1,\dots,N_l}$. The goal is to find a transforma-
 195 tion between the spaces of exact and proxy responses. (Notice that the *varimax* rotation
 196 does not affect the representation of the curves, but might affect the quality of the trans-
 197 formation).

Here, we restrict ourselves to functional linear regression models of the form given in Eq. 1. Training such a functional linear model in full generality is not straightforward. A simple choice to restrict the class of linear regression models is to postulate that, at any time t , $\tilde{y}_i(t)$ depends on $\tilde{x}_i(t)$ solely through its value at that time t . This assumption

leads to the Concurrent model

$$\tilde{y}_i(t) = \beta_0(t) + \tilde{x}_i(t)\beta_i(t) + \varepsilon_i(t), \tag{13}$$

198 which is a particular case of the functional linear model in Eq. 1 and corresponds to
 199 $T(x_i)(t) = \beta_0(t) + x_i(t)\beta_i(t)$. The Concurrent model will be used as baseline in our nu-
 200 merical application, and compared to our FPCA-based prediction approach.

201
 202 To simplify the exposition, in the following we assume that the same number of har-
 203 monics is retained for the two spaces, i.e., $D = D_{ex} = D_{app}$. However, the number of
 204 harmonics depends on the inherent variability of the learning set, which can be different
 205 for the exact and proxy responses. Ultimately, the number of harmonics to be employed
 206 depends on how rapidly the eigenvalues of the FPCA decomposition decrease for the spe-
 207 cific problem. It has to be chosen large enough to guarantee an exhaustive representation
 208 of the variability of the response curves, but small enough with respect to the number of
 209 elements in the learning set to avoid over-fitting when the regression model is constructed.

Given $N_l \leq N_r$ pairs of accurate and proxy responses, $\{(\tilde{x}_i(t), \tilde{y}_i(t))\}_{i=1, \dots, N_l}$, we postulate that there exists a $(D + 1) \times D$ matrix of real-valued coefficients $\boldsymbol{\beta}$ (with line index starting at 0, by convention) and a $N_l \times D$ error matrix \mathbf{E} , such that for any $(i, j) \in [1, N_l] \times [1, D]$,

$$c_{ij} = \beta_{0j} + \sum_{\ell=1}^D b_{i\ell}\beta_{\ell j} + e_{ij}, \tag{14}$$

where β_{ij} and e_{ij} are the components of $\boldsymbol{\beta}$ and \mathbf{E} , respectively. The errors, e_{ij} , are implicitly assumed to be Gaussian with zero mean and variance σ_j^2 , which depends only

on j . In matrix notation, the statistical model reads

$$\mathbf{C} = \mathbf{B}\boldsymbol{\beta} + \mathbf{E}, \quad (15)$$

210 where \mathbf{C} is the $N_l \times D$ matrix containing the scores of the exact responses, c_{ij} , and \mathbf{B} is
 211 the $N_l \times (D + 1)$ with elements of the first column $b_{i0} = 1$ by convention, and containing
 212 the scores of the proxy responses $b_{i(j-1)}$.

In the statistics literature, solving Eq. 15 for the coefficient matrix $\boldsymbol{\beta}$ is referred to as a
multivariate multiple regression problem ([Fox and Weisberg, 2011; Hastie et al., 2009]).

A simpler regression problem can be obtained by separating the regression models for the
 D responses, hence solving D independent regression problems

$$\mathbf{C}_{(j)} = \mathbf{B}\boldsymbol{\beta}'_{(j)} + \mathbf{E}'_{(j)} \quad (1 \leq j \leq D), \quad (16)$$

where $\mathbf{C}_{(j)}$ is the j^{th} column of the score matrix \mathbf{C} . A very convenient fact is that
 the columns of the Ordinary Least Squares (OLS) estimator of $\boldsymbol{\beta}$ coincides with the
 concatenated OLS estimators of $\boldsymbol{\beta}'_{(j)}$ [Hastie et al., 2009], that is

$$\widehat{\boldsymbol{\beta}}_{(j)} = \widehat{\boldsymbol{\beta}}'_{(j)} \quad (1 \leq j \leq D), \quad (17)$$

213 where $\widehat{\boldsymbol{\beta}}_{(j)}$ are the columns of the OLS estimator $\widehat{\boldsymbol{\beta}}$ (hereafter, the hat denotes the
 214 OLS estimator of the quantity). However, test statistics and confidence bands of the
 215 multivariate regression model cannot be directly derived from those obtained for the
 216 multiple linear regressions in Eq. 16 and have to be computed for the general regression
 217 model in Eq. 15. The formula of the simultaneous confidence bands is given in appendix
 218 C, together with a brief outline of the derivation.

3.4. Prediction of the exact response from the proxy response

Once the OLS estimator $\hat{\beta}$ has been obtained, the regression model is used to predict the exact response for all N_r geostatistical realizations on the basis of the corresponding proxy responses $\tilde{x}_i(t)$. The predicted exact response for the i^{th} realization is

$$\hat{y}_i(t) = \bar{y}(t) + \sum_{j=1}^D \hat{c}_{ij} \eta_j(t). \tag{18}$$

where

$$\hat{c}_{ij} = \hat{\beta}_{0j} + \sum_{\ell=1}^D \hat{\beta}_{j\ell} b_{i\ell}, \tag{19}$$

are the estimates of the scores with respect to the rotated harmonics.

The estimator of the linear regression model allows us to predict the \hat{c}_{ij} scores solely from the scores b_{ij} of the proxy responses, hence predicting $\tilde{y}_i(t)$ without solving the exact model. We emphasize the difference between the proxy response $x_i(t)$ (or $\tilde{x}_i(t)$, which is the projection onto the lower dimensional space defined by the first D harmonics, $\{\zeta_j\}_{j=1,\dots,D}$), and the predicted exact response $\hat{y}_i(t)$: they both approximate the “true” response $y_i(t)$, but, while $x_i(t)$ is simply the result of the proxy model and lives in the space defined by the basis of the proxy curves, $\hat{y}_i(t)$ results from applying the error models to the proxy response and lives in the space of the exact responses (more precisely: in the subspace defined by the orthonormal basis formed by the first D harmonics, $\{\eta_j\}_{j=1,\dots,D}$).

Surrogating $y(t)$ by $\hat{y}(t)$ is prone to errors: first, $\{\eta_i(t)\}_{i=1,\dots,N_i}$ depends on the quality of the learning set; second, the subspace of the prediction is further reduced by considering only the first D harmonics; third, the coefficients \hat{c}_{ij} are predicted through the OLS estimator of a linear regression model, and thus entails statistical uncertainties and possibly systematic errors due to the choice of a simple linear model.

4. Numerical test case: An idealized NAPL pollution problem

234 We consider an idealized groundwater pollution problem in which the fate of a NAPL
235 plume has to be predicted. We model a portion of aquifer as a vertical 2D domain of
236 length $10.8m$ and depth $5.1m$ discretized into cells of size $10cm \times 10cm$. Gravity effects are
237 neglected, which implies that the density of the NAPL phase is equal to the water density.
238 No-flow boundary conditions are imposed at the upper and lower boundaries, whereas the
239 pressure is fixed at the right boundary. The contaminant is released at the left boundary
240 (a constant influx is assigned) and displaces the water initially present in aquifer. We are
241 interested in the time evolution of NAPL saturation at the right boundary. Two cases are
242 investigated; first, we estimate the uncertainty on the contaminant breakthrough curve
243 computed by averaging the saturation along the right boundary; then, we consider a
244 single-point breakthrough curve obtained by sampling the saturation in a single cell (Sec.
245 4.5.2). As the NAPL is immiscible with water, the exact model solves the multiphase flow
246 and transport equations, which require solving a pressure equation and a highly nonlinear
247 phase-transport equation [see, e.g., *Marle*, 1981; *Helmig*, 1997]. The two equations are
248 highly coupled and characterized by fluxes that exhibit a non-linear dependence on NAPL
249 saturation. (The full system of equations is described in Appendix A.)

250 The uncertainty on the transport properties of the aquifer (permeability and porosity)
251 is represented by a set of $N_r = 1000$ geostatistical realizations that are generated by a
252 multipoint geostatistical method (DeeSse) [*Mariethoz et al.*, 2010] with a training image
253 obtained from data of facies-distribution collected at the Herten site (Germany) [*Bayer*
254 *et al.*, 2011]. As an example, three realizations are shown in Fig. 3.

4.1. The proxy model

255 The proxy model simplifies the physics of the problem by treating the NAPL as an ideal
256 tracer, thus solving a linear transport problem. Although it is possible to further improve
257 the computational efficiency by simplifying the description of the heterogeneity (e.g., by
258 some upscaling or multiscale methods [see, e.g., *Josset and Lunati, 2013*]), here we do not
259 approximate the aquifer properties.

260 In practical situations, replacing a multiphase flow problem by a single-phase (tracer-
261 transport) problem considerably reduces the computational costs. Indeed, a large part of
262 the cost of solving the flow and transport system stems from the solution of the elliptic
263 (or parabolic) equation that governs the pressure. Due to the effects of the saturation on
264 the fluxes, this equation has to be solved at every time step in multiphase problems. In
265 contrast, if the pollutant is considered as an ideal tracer, the saturation does not impact
266 the velocity, and the pressure equation has to be solved only once. The NAPL transport
267 equation becomes linear and can be solved very efficiently by streamline methods (here, we
268 use a Finite-Volume upwind scheme that can be seen, in some sense, as a very rudimentary
269 streamline method without sub-grid interpolation of the velocity field).

4.2. The learning dataset

270 After the proxy responses have been obtained by solving the ideal transport problem
271 and computing the contaminant breakthrough curves for the whole sample of 1000 real-
272 izations, we construct the learning set by identifying a subset of $N_l = 20$ realizations. The
273 realizations can be selected in several ways, including a simple random choice. Here, we
274 use a clustering technique to group the proxy responses based on their l_2 -distance, and
275 we choose the k -medoid curves as representative of the clusters (Distance Kernel Method

276 [*Scheidt and Caers, 2009a*]). The medoids define the subset of realizations, $\{R_i\}_{i=1,\dots,N_l=20}$,
277 for which the exact responses are computed by solving the multiphase transport problem.
278 Additional tests (not reported here) with learning sets consisting of $N_l = 50$ and $N_l = 100$
279 realizations did not show a significant improvement of the quality of the learning set. This
280 suggests that only 20 realizations are sufficient to obtain a satisfactory error model for
281 the present test case. Cross validation tests can be performed to identify the optimal size
282 of the learning set.

283 As the numerical NAPL breakthrough curves are discrete in time, a spline basis is
284 defined to interpolate the discrete data and construct the functional objects. In the present
285 test case, data points are fairly smooth and a rather small number of basis functions is
286 necessary for an accurate representation of the data (here, only 50 splines are used as basis
287 functions). The 20 pairs of spline-interpolated proxy and exact curves in the learning set,
288 $\{(x_i(t), y_i(t))\}_{i=1,\dots,N_l=20}$, are shown in Fig. 4.

4.3. Understanding the data using FPCA

289 To extract the relevant information from the data and to reduce the problem dimension-
290 ality, we apply FPCA independently to both sets of approximate and exact curves in the
291 learning set. As in standard PCA, if all the components (harmonics) are considered, no
292 approximation is made and the data are represented exactly. However, the eigenvalues of
293 higher order harmonics decrease so fast that the first three components describe more than
294 97% and 99% of the variability of proxy and exact curves, respectively. In the subspaces
295 defined by the first three harmonics, each curve is described by the corresponding three
296 scores and by the sample means. To improve the interpretability of the data, a rotation

297 is sought with the varimax algorithm [*Ramsay et al.*, 2012]. The rotated harmonics for
 298 both sets of curves are shown in Fig. 5.

299 In the subset of the exact responses, the first rotated component explains the devi-
 300 ation from the mean behavior measured at late time. The second rotated component
 301 describes the variation at the beginning of the breakthrough curve, thus enlightening high-
 302 connectivity paths. The third component explains the variation observed at intermediate
 303 time. In the proxy subset, the first rotated component describes the initial variability; the
 304 second component highlights the variation at high saturation; and the third component
 305 explains the variation observed at intermediate time. By analyzing the projection of the
 306 curves on these components, it is possible to gain information about the data, for instance
 307 about the link between the early-time responses and the late-time variations. We refer to
 308 *Henderson* [2006] for an example in hydrology.

4.4. Regression model and evaluation of the proxy

309 The linear regression model is built between the scores of proxy and exact curves, which
 310 represent their coordinates with respect to the two orthonormal bases formed by the first
 311 three harmonics. Three linear regression problems (one for each exact-response score,
 312 $j = 1, 2, 3$) are solved to establish a relationship with the three proxy-response scores.
 313 The resulting coefficients of the three regression models are

	β_{0j}	β_{1j}	β_{2j}	β_{3j}	R^2	p-value
$j = 1$	$-2.3 \cdot 10^{-16}$	0.42	0.18	-0.37	0.99	$< 2 \cdot 10^{-16}$
$j = 2$	$4.4 \cdot 10^{-17}$	0.82	-0.02	0.37	0.99	$< 2 \cdot 10^{-16}$
$j = 3$	$1.6 \cdot 10^{-16}$	0.51	0.03	0.08	0.97	$1.3 \cdot 10^{-12}$

(20)

314 Notice that the R^2 values are quite high and that $\beta_{0j} \approx 0$, which suggests that the linear
 315 regression model preserves the mean. The dependency among scores is illustrated in Fig.

316 6. The relationships between the scores of the three harmonics of the exact curves and the
317 scores of the first harmonic of the proxy curves are rather well approximated by the linear
318 regression. The scores of the second harmonic of the proxy curves are less important as
319 it is indicated by the low values of β_{22} and β_{23} . This might be due to the fact that the
320 proxy second harmonic explains the variability of the curves for saturations close to one,
321 a situation that is not observed in the two-phase responses.

4.5. Performance of the regression model as error model

322 In general, the proxy-curve scores are informative of the exact-curve scores, at least
323 for the curves pairs in the learning set. This suggests that, despite the rather primitive
324 physical model employed, the regression model can be effectively used to predict the
325 exact responses of the realizations for which only the proxy solution is available. The
326 exact response is predicted on Eqs. 18 and 19.

4.5.1. Prediction of the average breakthrough curve at the outlet

328 We start by considering the prediction of the breakthrough curve calculated by averaging
329 the saturation at the right-hand boundary. Examples of two predicted curves are shown in
330 Fig. 7a and b. Despite the fact that the curves are very different for the two realizations,
331 both predictions are in good agreement with the exact responses. In general, the behaviour
332 of the exact response is well predicted, with the exception of some fluctuations at early
333 times. The error model greatly improves the proxy solution and provides a much better
334 prediction than the Concurrent model, which is unable to significantly modify the shape
335 of the curves due to the use of only concurrent information.

336 The differences between predicted and exact curves are illustrated in Fig. 7c for all
337 $N_r = 1000$ realizations, together with the mean error. The maximum differences in the

338 saturation are observed at early time and are about 10%; later, the saturation discrepancy
339 remains below $\pm 1.8\%$ for 68% of the realizations and below $\pm 4\%$ in the worst cases. The
340 mean error is very close to zero, which shows that the predicted curves conserve the
341 mean behaviour of the exact curves, and that the subset of 20 realizations selected in the
342 learning set is representative of the whole sample to describe the mean behaviour.

343 Fig. 8a shows the histograms for the l_2 -norm and the l_∞ -norm of the errors. We compare
344 the performance of the error model based on FPCA with the Concurrent functional linear
345 regression model. The histogram of the l_∞ -norm shows that on average the maximum
346 deviation is 4.5% for FPCA, and about 8% for the Concurrent model. The l_2 -error is on
347 average more than three times lower for the FPCA-based model.

348 In many applications, the uncertainty is quantified in terms of the quantiles of the
349 responses. Fig. 7d displays the quantile curves obtained using the different models. The
350 Concurrent model fails to reproduce the 90th percentiles, because it is unable to modify
351 the plateau of the proxy curves close to saturation one; it performs better for the other
352 quantiles. The quantiles curves computed using only the learning set of exact responses
353 (as suggested by [*Scheidt and Caers, 2009a, b*]), are slightly biased estimates of the exact
354 quantiles. An excellent estimate is obtained with the functional error model, which is able
355 to correct the approximate responses and predicts quantiles close to the exact ones.

356 4.5.2. Prediction of single-point breakthrough curve

357 In this second test case, we are interested in predicting the breakthrough curve of the
358 contaminant at a precise location, defined by a single cell of the numerical grid, which is
359 located at mid-depth at the outlet. In contrast to the breakthrough curves averaged over
360 the whole outlet, in which the effects of extreme permeability structures (flow barriers

361 or preferential pathways) are smoothed, the single-point breakthrough curves display a
362 variety of shapes. The large contrast in permeability and in connectivity at the sampling
363 location leads to important differences, particularly in the first arrival time.

364 In this case, it is useful to apply a translation in time to redefine the origin, which is
365 chosen to be the first arrival time. This procedure is referred to as registration in the
366 FDA literature *Ramsay* [2006]; *Ramsay et al.* [2009]. For the translated responses in the
367 learning set FPCA is then applied and the dimensionality is reduced as described above.
368 Again, we use the first three harmonics, which describe more than 98% of the variability
369 of the shape of the curves after the registration. An example of proxy, predicted and exact
370 curves after registration is shown in Fig. 9a for a realization that does not belong to the
371 learning set.

372 Beside the prediction of the shape, it is now necessary to predict the first arrival time
373 and translate back the predicted curves. The first arrival time is predicted jointly to the
374 scores of the harmonics by solving a 4×4 regression model, where the 4th dimension is
375 the first arrival times of the proxy responses, which have been used for the registration.
376 Fig. 9b compares the proxy and exact curves with the predicted curve after translation
377 by the predicted arrival time (these curves correspond to the registered curves in Fig. 9a).
378 For the whole sample of realizations, the mean saturation error is close to zero and with
379 a standard deviation that remains below 0.04 (Fig. 9c).

380 The predicted quantile curves (shown in Fig. 9d) are in good agreement with the exact
381 quantile curves for P50 and P90, but P10 is biased. As the concurrent model would
382 perform very poorly in this case because it is unable to deal with curves characterized by
383 different arrival times, we compare our methodology with the quantile curves obtained

384 directly from the exact response in the learning set (this procedure corresponds to the
385 classical DKM). As both the functional error model and the DKM estimates depend on
386 the clustering, we have applied both methodologies 200 times. The example shown in Fig.
387 9 is representative of the typical behaviors of the methods (i.e., the quantiles are close to
388 the average quantiles obtained from the 200 applications of the methods shown in Fig. 9e
389 and d. In average, the functional error model is more robust than DKM and provides a
390 better prediction of the P10 quantile curve.

391 4.5.3. Effects of the number of harmonics

392 Here, we investigate the effects of the number of harmonics on the prediction of single-
393 point breakthrough curves. In order to increase the difficulty of the problem, we do
394 not apply the registration as in the previous section (i.e., the breakthrough curves are
395 not translated by their first arrival times). On one hand this requires more harmonics
396 to describe the variability of the curves; on the other hand it allows us to demonstrate
397 that the functional error model is able to correct for different arrival times also without
398 registration.

399 We consider 200 different learning sets, which are selected by DKM clustering with
400 different initialization. For each learning set we apply FPCA and then construct the
401 functional error models by employing a different number of harmonics. The quality of the
402 prediction is measured by the l_2 distance between the predicted and exact responses for
403 all 1000 realizations.

404 The performance of the method (expressed as median error and confidence interval
405 of the responses of the 200 learning sets) is presented in Fig. 10 as a function of the
406 number of harmonics. The error exhibits a minimum around 5-7 harmonics. Indeed,

407 when the number of harmonics is increased from 2 to 5, the variability of the learning set
408 represented increases from 92% to 99%, leading to an improved error model. If the number
409 of harmonics is increased further, the error increases quite rapidly. For 12 harmonics
410 errors are very large and fluctuate greatly depending on the choice of the learning set.
411 This behavior is a clear signature of over-fitting, as the large number of harmonics is
412 not balanced by the size of the learning set (consisting of 20 pairs of curves) and the
413 parameters of the regression model are not constrained enough by the data.

5. Conclusions

414 We have presented a novel methodology that combines elements of Functional Data
415 Analysis and Machine Learning to construct error models that improve uncertainty quan-
416 tification. The approach is purpose-oriented as it is formulated directly on the quantity
417 of interest (in the case considered here, the contaminant breakthrough curve) rather than
418 on the state of the system (e.g., the entire saturation and pressure fields).

419 The core idea of the method is to construct an error model from a learning set containing
420 pairs of proxy and exact responses of a subset of realizations, and to predict the exact
421 responses of the entire sample without solving the exact model for all realizations. FPCA
422 is employed to separately reduce the dimensionality of the spaces of exact and proxy
423 responses in the learning set. The advantage is twofold: on one hand, the small dimension
424 allows a diagnostic of the regression model on scores to assess the informativeness of the
425 proxy for the application at hand; on the other hand, using spaces of lower dimension
426 reduces the risk of over-fitting when the regression model is constructed.

427 The method has been tested for a synthetic contamination problem, in which the break-
428 through curve of a NAPL contaminant is predicted with the help of a tracer transport

429 simulation (as proxy model). We have obtained excellent results with a learning set con-
430 sisting of 20 pairs of curves (corresponding to 20 realizations out of a sample of 1000) and
431 considering only the first three harmonics, which describe more than 97% of the variabil-
432 ity. Visual inspection of the score scatter plots shows that the proxy is indeed potentially
433 very informative of the exact response (this is confirmed by a linear determination coef-
434 ficient $R^2 = 0.97$). Notice that this is not necessarily an indication of the quality of the
435 predictions as the size of the learning set and the number of harmonics also influence the
436 accuracy of the prediction. For both test case, the error model allows us to solve a two-
437 phase problem only for the 20 realizations, whereas a simple tracer transport problem is
438 solved for all realizations in the sample. The gain in computational efficiency is evident as
439 multiphase transport requires solving the pressure problem at every time step, in contrast
440 to ideal tracer transport, which requires solving the pressure equation only once.

441 In comparison to the Concurrent model (an existing methodology used to correct proxy
442 responses), we have demonstrated an error reduction by a factor 3 when the functional
443 error model is employed. Also, the error model improves the uncertainty quantification
444 with respect to the estimate obtained solely on the basis of the 20 exact responses in the
445 learning set (this approach corresponds to the DKM, which uses the proxy responses only
446 to cluster the realizations). Beside an increase in accuracy, the methodology presents two
447 advantages over the DKM. First, the error model allows us to use the proxy response to
448 predict the exact response for any new geostatistical realization that might be successively
449 generated; this clearly opens new possibilities to use the model beyond the context of
450 uncertainty quantification, and in particular for Bayesian inference, model calibration
451 and optimization. Second, simultaneous confidence bands of the predicted curves can be

452 defined by propagating the errors of the multivariate regression model. Notice that the
453 residual uncertainty due to the size of the learning set and to the truncation of the basis
454 should be taken into account.

455 Combining FPCA and machine learning can be seen as a general framework in which
456 each component can be modified and improved, if it is required to improve accuracy. For
457 instance, the rather crude linear regression model between the three-dimensional spaces
458 of exact and proxy responses can be made more complex by increasing the dimensions
459 (possibly with different truncations for the proxy and the exact model) or by refining
460 the mathematical form of the statistical model to predict the scores. Possible enhance-
461 ments include linear regression models with more complex basis functions (polynomials or
462 others), but may also entail kernel methods like co-kriging. Almost any multivariate pre-
463 diction may be adapted to this problem once the dimensionality reduction is performed.
464 Another potential improvement is to perform the dimensionality reduction jointly for the
465 proxy and the exact spaces, in order to optimize the informativeness of the proxy rather
466 than the description of the variability of each response space independently. Indeed, in
467 very complex test cases, it might occur that some small-eigenvalue harmonics of the proxy
468 response might explain large-eigenvalue characteristics of the exact curves. This can be
469 done by replacing FPCA by Functional Canonical Correlation Analysis [*Ramsay*, 2006]
470 or by Functional Partial Least Squares [*Cuevas*, 2014].

471 Finally, we observe that the proposed framework can be applied far beyond the con-
472 tamination example that we have presented. It can be useful in virtually any situation
473 in which the most reliable technique has to be surrogated by an approximate method.
474 Applications are not limited to the case in which evaluating exact response involves the

475 solution of a complex numerical model, but also to situations in which the proxy or the
476 exact responses consist of experimental data. The FDA framework would be then impor-
477 tant to compare information with different temporal resolutions. Also, the error model
478 can potentially be very useful in the context of Bayesian inference, when the number of
479 responses that have to be evaluated (e.g., in Metropolis-Hastings algorithms and alike)
480 is typically of the order of 10^5 . In this case, a functional error model capable to predict
481 the exact responses only on the basis of the proxy responses can substantially speed up
482 MCMC algorithms, as it reduces the cost of likelihood estimation. This would improve
483 the efficiency of the calibration and optimization algorithms, which are often used in
484 hydrogeological applications.

Appendix A: Multiphase and single-phase transport equations

485 Assuming that both phases are incompressible and neglecting gravity and capillary
486 effects, the saturation of the NAPL, S , is governed by the following system of equations:

$$\nabla \cdot \left[\left(\frac{k_n(S)}{\mu_n} + \frac{k_w(1-S)}{\mu_w} \right) k \nabla p \right] = 0, \quad (\text{A1})$$

$$\frac{\partial}{\partial t}(\phi S) - \nabla \cdot \left(\frac{k_n(S)}{\mu_n} k \nabla p \right) = 0, \quad (\text{A2})$$

487 where the absolute permeability, k , and the porosity, ϕ , are aquifer properties; p is the
488 pressure; μ_n and μ_w are the viscosities of NAPL and water, respectively; and k_n and k_w are
489 the relative permeabilities of NAPL and water, respectively, which are nonlinear functions
490 of the saturation. Together with the constitutive relationships for the permeabilities (here,
491 they are assumed quadratic i.e., $k_n(S) = S^2$ and $k_w(S) = (1 - S)^2$), the two equations
492 above form a complete system of equations that can be solved for p and S to calculate
493 the NAPL breakthrough curves. These curves are the responses of the exact (multiphase)
494 model.

495 Due to the nonlinearity of the relative permeability, the system above is computationally
496 expensive because the two equations are coupled and the pressure equation has to be
497 solved at any time step. This problem can be avoided by neglecting the nonlinearity of
498 the permeabilities, hence approximating the system above as

$$\nabla \cdot \left(\frac{k}{\mu_w} \nabla p \right) = 0, \quad (\text{A3})$$

$$\frac{\partial}{\partial t}(\phi S) - \nabla \cdot \left(S \frac{k}{\mu_w} \nabla p \right) = 0, \quad (\text{A4})$$

499 which corresponds to a simple tracer transport problem without mechanical dispersion.

Appendix B: Linear models for functional responses with functional predictors

A simple class of linear models is the Concurrent model [Ramsay, 2006]. The value of the response variable $y(t)$ is predicted solely by the value of the functional covariate at the same time t

$$y_i(t) = \alpha(t) + x_i(t)\beta(t) + \varepsilon_i(t), \quad (\text{B1})$$

500 where $\varepsilon_i(t)$ are the functional errors and the functions $\alpha(t)$ and $\beta(t)$ are estimated by
 501 minimizing the sum of squares under some penalty on the roughness of the functions to
 502 avoid overfitting and loose predictability power. Despite the rather arbitrary choice of
 503 the degree of smoothness of the functional parameters, this method is quite fast but also
 504 rudimentary because there is a priori no reason to assume that only concurrent features
 505 of the curves are relevant (this is well illustrated by the synthetic test to predict the
 506 single-point breakthrough curve in Sec. 4.5.2).

507

A generalized formulation is when the functional variable contributes to the prediction for all possible time values

$$y_i(t) = \alpha(t) + \int x_i(s)\beta(s, t)ds + \varepsilon_i(t) \quad (\text{B2})$$

508 which allows the predicted response to depend on the functional covariate at all times,
 509 but $\beta(s, t)$ is now bivariate. The application of this model is known to be particularly
 510 challenging as the smoothing constraints to be imposed is of paramount importance.

Appendix C: Simultaneous confidence bands for multiple multivariate linear regression

To take into account the uncertainty stemming from the linear regression, we derive simultaneous confidence bands for the predicted curve $\hat{y} = \mathbf{b}'\hat{\beta}\boldsymbol{\eta}(t)$, where $1 - \alpha$ is the level of confidence that the exact curve $\tilde{y}(t) = \mathbf{b}'\beta\boldsymbol{\eta}(t)$ is within the confidence bands for all t , that is

$$\Pr\left(\tilde{y}(t) \in [\hat{y}(t) - w_\alpha(t), \hat{y}(t) + w_\alpha(t)] \text{ for all } t\right) = 1 - \alpha \quad (\text{C1})$$

511 and, following the sketch of proof below, where $D_{ex} + D_{app} < N_l$ is assumed,

$$\begin{aligned} w_\alpha(t) &= \sqrt{\frac{D_{ex}(N_l - D_{app} - 1)}{N_l - D_{ex} - D_{app}}} F_{D_{ex}, N_l - D_{ex} - D_{app}}(\alpha) \\ &\times \sqrt{(1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b}) \frac{N_l}{N_l - D_{app} - 1} \boldsymbol{\eta}'(t) \hat{\boldsymbol{\Sigma}} \boldsymbol{\eta}(t)}, \end{aligned} \quad (\text{C2})$$

512 where $\boldsymbol{\eta}(t)$ the values of the exact harmonics; $F(\alpha)$ Fisher's α -quantile; and $\hat{\boldsymbol{\Sigma}}$ the
513 covariance matrix of the errors estimated on the learning set.

The key step of the derivation is the use of Scheffe's Lemma that states that, for a symmetric and positive definite matrix $\Gamma \in \mathbb{R}^{p \times p}$, the following statements are equivalent for any vector $\mathbf{v} \in \mathbb{R}^p$ and constant $c > 0$

$$\left(\mathbf{v}'\Gamma\mathbf{v} \leq c^2\right) \iff \left(|\boldsymbol{\psi}'\mathbf{v}| \leq c\sqrt{\boldsymbol{\psi}'\Gamma^{-1}\boldsymbol{\psi}} \quad \forall \boldsymbol{\psi} \in \mathbb{R}^p\right) \quad (\text{C3})$$

Sketch of proof

The residuals $\hat{\mathbf{E}} = \hat{\mathbf{C}} - \mathbf{C}$ are centred and with covariance $\mathbb{E}[\hat{\mathbf{E}}'\hat{\mathbf{E}}] = (N_l - D_{app} - 1)\boldsymbol{\Sigma}$, where $(\boldsymbol{\Sigma})_{jk} = \sigma_{jk}$. Assuming that \mathbf{E} is Gaussian entails that $\hat{\boldsymbol{\beta}}$ is Gaussian, whereof $\mathbf{c} \sim \mathcal{N}_{D_{ex}}\left(\mathbf{b}'\beta, (1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b})\boldsymbol{\Sigma}\right)$. Then $\left(\frac{\mathbf{b}'\hat{\beta} - \mathbf{b}'\beta}{\sqrt{1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b}}}\right)' \left(\frac{1}{N_l - D_{app} - 1}\boldsymbol{\Sigma}\right)^{-1} \left(\frac{\mathbf{b}'\hat{\beta} - \mathbf{b}'\beta}{\sqrt{1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b}}}\right)'$ follows a Chi-squared distribution $\chi_{D_{ex}}^2$. On the other hand, the usual estimator $\hat{\boldsymbol{\Sigma}}$ of $\boldsymbol{\Sigma}$ follows a Wishart distribution independently from $\hat{\beta}$. We then obtain the following

$$t^2 = \left(\frac{\mathbf{b}'\hat{\beta} - \mathbf{b}'\beta}{\sqrt{1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b}}}\right)' \left(\frac{N_l}{N_l - D_{app} - 1}\hat{\boldsymbol{\Sigma}}\right)^{-1} \left(\frac{\mathbf{b}'\hat{\beta} - \mathbf{b}'\beta}{\sqrt{1 + \mathbf{b}'(\mathbf{B}'\mathbf{B})^{-1}\mathbf{b}}}\right) \sim T_{D_{ex}, N_l - D_{app} - 1}^2. \quad (\text{C4})$$

As the Hotelling T^2 -distribution can be expressed in term of the F -distribution, we can write that, with probability $1 - \alpha$,

$$t^2 \leq \frac{D_{ex}(N_l - D_{app} - 1)}{N_l - D_{ex} - D_{app}} F_{D_{ex}, N_l - D_{ex} - D_{app}}(\alpha), \quad (\text{C5})$$

514 where $F_{p,q}(\alpha)$ stands for the α -quantile of the Fisher-Snedecor distribution with parame-
515 ters p and q .

516 Using Scheffe's Lemma (eq. C3) for $\mathbf{v} = \mathbf{b}'\hat{\boldsymbol{\beta}}$ the vector of predicted scores and $\boldsymbol{\psi}$ the
517 vector of the exact harmonics values $\boldsymbol{\eta}(t)$, the second statement gives us the simultaneous
518 confidence bands on the prediction.

519 Acknowledgments.

520 The data to support this article result from numerical simulations (multiphase and
521 tracer transport) performed with the open source code MaFlot (*Matlab Flow and Trans-*
522 *port* [Künze and Lunati, 2012, 2013]). The data treatment is performed with the *fda*
523 package [Ramsay et al., 2012] implemented in R [R Core Team, 2013]. Upon request by
524 email, the authors would provide the simulated data and codes.

525 The authors thank Rouven Künze for his assistance with the flow simulations, Guillaume
526 Pirot and Philippe Renard for providing the realizations of the hydraulic conductivity, and
527 Céline Scheidt for sharing her DKM code. Many thanks are due to L. Dmbgen for his
528 teachings, to V. Demyanov and A.H. Elsheikh for many useful discussions and to the
529 reviewers for their suggestions and careful editing.

530 This project is supported by the Swiss National Science Foundation as a part of the
531 ENSEMBLE project (Sinergia Grant No. CRSI22-132249/1). David Ginsbourger ac-
532 knowledges support from the Institute of Mathematical Statistics and Actuarial Science,

533 University of Bern. Ivan Lunati is Swiss National Science Foundation (SNSF) Professor
534 at the University of Lausanne (SNSF grant numbers PP00P2-123419/1 and PP00P2-
535 144922/1).

References

- 536 Ballin, P., A. Journel, and K. Aziz (1992), Prediction of uncertainty in reservoir perfor-
537 mance forecast, *Journal of Canadian Petroleum Technology*, 31(4).
- 538 Bayer, P., P. Huggenberger, P. Renard, and A. Comunian (2011), Three-dimensional
539 high resolution fluvio-glacial aquifer analog: Part 1: Field study, *Journal of Hydrology*,
540 405(1), 1–9.
- 541 Christie, M. (1996), Upscaling for reservoir simulation, *Journal of Petroleum Technology*,
542 48(11), 1004–1010.
- 543 Christie, M. A., J. Glimm, J. W. Grove, D. M. Higdon, D. H. Sharp, and M. M. Wood-
544 Schultz (2005), Error analysis and simulations of complex phenomena, *Los Alamos*
545 *Science*, 29(6).
- 546 Cuevas, A., Febrero, M., and Fraiman, R. (2002). Linear functional regression: the case
547 of fixed design and functional response, *Canadian Journal of Statistics*, 30(2), 285-300.
- 548 Cuevas, A. (2014), A partial overview of the theory of statistics with functional data,
549 *Journal of Statistical Planning and Inference*, 147, pp. 1–23.
- 550 Dagan, G. (2002), An overview of stochastic modeling of groundwater flow and transport:
551 From theory to applications, *Eos, Transactions American Geophysical Union*, 83(53),
552 621.
- 553 Durlofsky, L. (2005), Upscaling and gridding of fine scale geological models for flow sim-
554 ulation, in *8th International Forum on Reservoir Simulation Iles Borromees, Stresa*,

- 555 *Italy*, pp. 20–24.
- 556 Fox, J., and Weisberg, S. (2011). *Multivariate Linear Models in R*.
- 557 Hastie, T., R. Tibshirani, and J. Friedman (2009), *The Elements of Statistical Learning*,
- 558 Springer
- 559 Helmig R. (1997) *Multiphase Flow and Transport Processes in the Subsurface*, Springer
- 560 Verlag, Berlin-Heidelberg
- 561 Henderson, B. (2006), Exploring between site differences in water quality trends: a func-
- 562 tional data analysis approach, *Environmetrics*, 17(1), 65–80.
- 563 Josset, L., and I. Lunati (2013), Local and global error models to improve uncertainty
- 564 quantification, *Mathematical Geosciences*, pp. 1–20.
- 565 Kaiser, H. (1958), The varimax criterion for analytic rotation in factor analysis, *Psy-*
- 566 *chometrika*, 23(3), 187–200, doi:10.1007/BF02289233.
- 567 Künze, R., and I. Lunati (2013), Adaptive multiscale simulations of density-driven insta-
- 568 bilities, *Journal of Computational Physics*, 255, 502–520, 10.1016/j.jcp.2012.02.025.
- 569 Künze, R., and I. Lunati (2012), MaFloT - Matlab Flow and Transport. Published under
- 570 the GNU license agreement on www.mafлот.com
- 571 Mariethoz, G., P. Renard, and J. Straubhaar (2010), The direct sampling method to
- 572 perform multiple-point geostatistical simulations, *Water Resources Research*, 46(11),
- 573 W11,536.
- 574 Marle, C.M. (1981), *Multiphase flow in porous media*, Institut Francais du Petrole, Paris,
- 575 France.
- 576 McLennan, J., and C. Deutsch (2005), Ranking geostatistical realizations by measures
- 577 of connectivity, in *SPE International Thermal Operations and Heavy Oil Symposium*,

- 578 98168, Alberta, Canada.
- 579 O’Sullivan, A., and M. Christie (2005), Error models for reducing history match bias,
580 *Computational Geosciences*, 9(2-3), 125–153.
- 581 O’Sullivan, A., and M. Christie (2006), Simulation error models for improved reservoir
582 prediction, *Reliability Engineering & System Safety*, 91(10), 1382–1389.
- 583 R Core Team (2013), *R: A Language and Environment for Statistical Computing*, R Foun-
584 dation for Statistical Computing, Vienna, Austria, ISBN 3-900051-07-0.
- 585 Ramsay, J. J. O., G. Hooker, and S. Graves (2009), *Functional data analysis with R and*
586 *MATLAB*, Springer.
- 587 Ramsay, J. O. (2006), *Functional data analysis*, Wiley Online Library.
- 588 Ramsay, J. O., H. Wickham, S. Graves, and G. Hooker (2012), *fda: Functional Data*
589 *Analysis*, r package version 2.3.2.
- 590 Renard, P., and G. de Marsily (1997), Calculating equivalent permeability: a review,
591 *Adv. Water Res.*, 20(5-6), 253–278.
- 592 Richman, M. B. (1986), Rotation of principal components, *Journal of climatology*, 6(3),
593 293–335.
- 594 Scheidt, C., and J. Caers (2009a), Representing spatial uncertainty using distances and
595 kernels, *Mathematical Geosciences*, 41(4), 397–419.
- 596 Scheidt, C., and J. Caers (2009b), Uncertainty quantification in reservoir performance
597 using distances and kernel methods—application to a west africa deepwater turbidite
598 reservoir, *SPE Journal*, 14(4), 680–692.
- 599 Ullah, S., C. F. Finch, et al. (2013), Applications of functional data analysis: A systematic
600 review, *BMC medical research methodology*, 13(1), 43.

601 Wen, X., and J. Gómez-Hernández (1996), Upscaling hydraulic conductivities in hetero-
602 geneous media: An overview, *J. Hydr.*, 183(1-2), R9–R32.

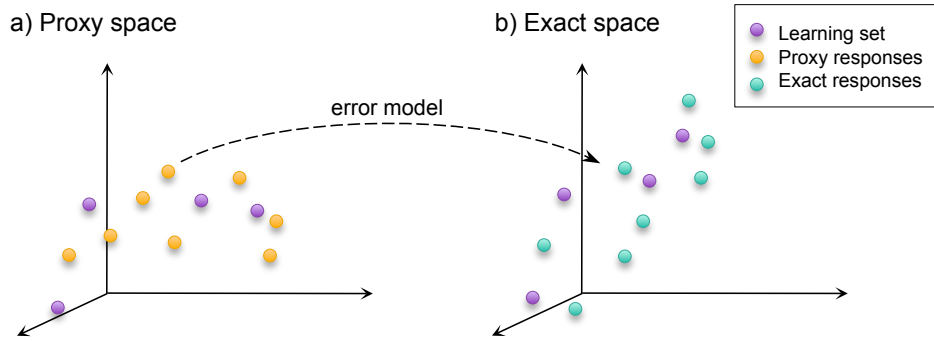


Figure 1. A statistical model is built on the learning set to relate the coefficients of the elements $x_i(t)$ in the proxy-response space to the coefficients of the elements $y_i(t)$ in the exact-response space. It is used as error model to predict the exact response from the proxy response.

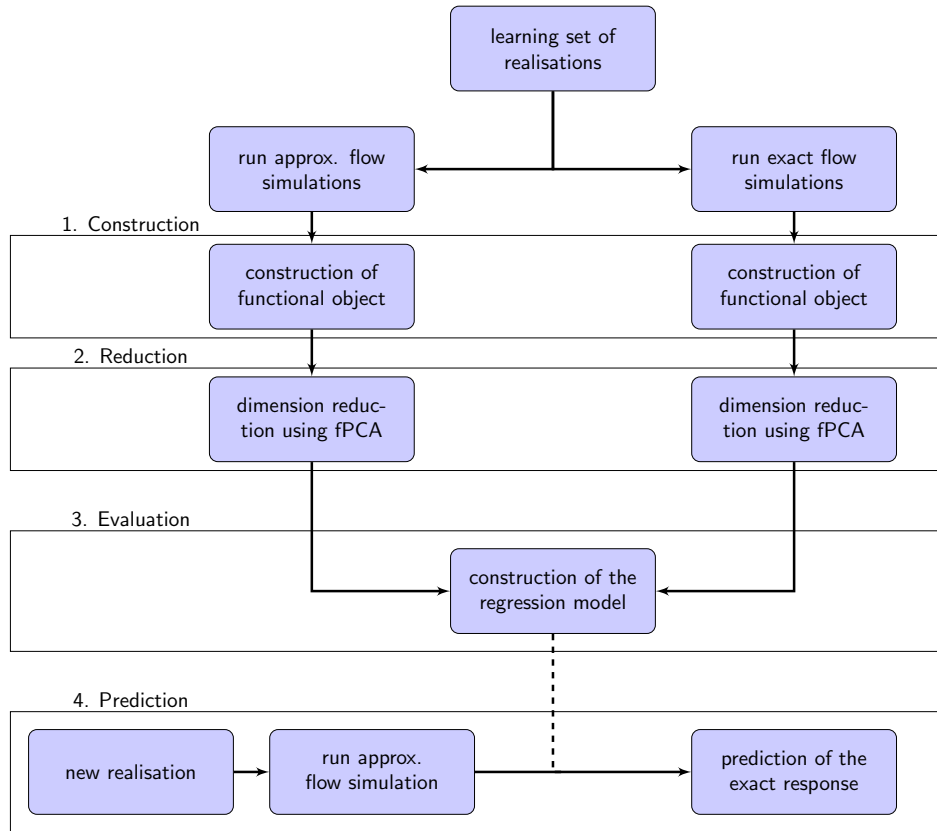


Figure 2. Flowchart of the methodology. After a learning set of realization has been constructed by selecting a subset of realizations and calculating pairs of proxy- and exact-response curves, the exact responses for the realizations that are not in the learning set can be predicted in four steps: 1. first, the functional objects are constructed by spline interpolation, 2. then, the dimensions of the subspaces of exact and proxy responses are reduced by means of FPCA, 3. next, a regression model is constructed between the proxy and the exact scores; 4. finally, the regression model is used to predict the exact responses of the realizations that are not in the learning set.

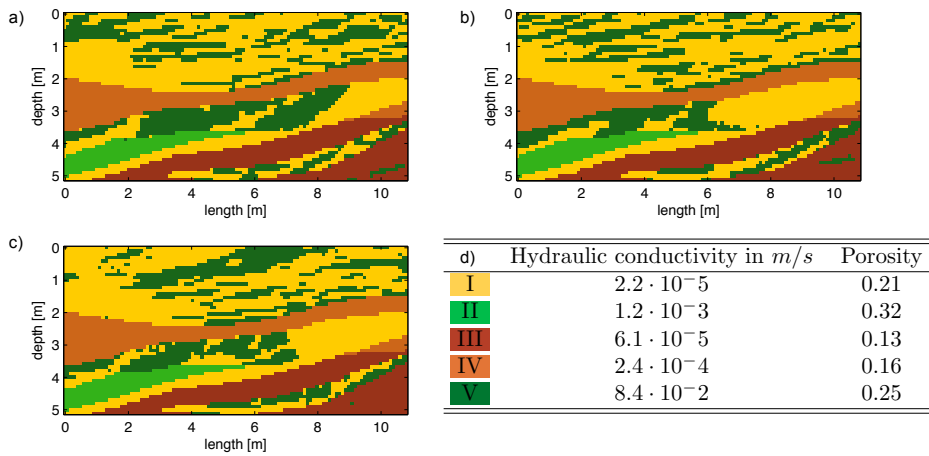


Figure 3. (a), (b) and (c): three examples of geostatistical realizations generated by a multipoint methods (DeeSSe, [Mariethoz *et al.*, 2010]) with training image from the Herten site (Germany) [Bayer *et al.*, 2011]. The different colors correspond to 5 different facies, whose properties are reported in (d). The three realizations belong to the set of realizations used to construct the learning set; the corresponding NAPL breakthrough curves obtained with the exact and with the approximate models are highlighted in Fig. 4.

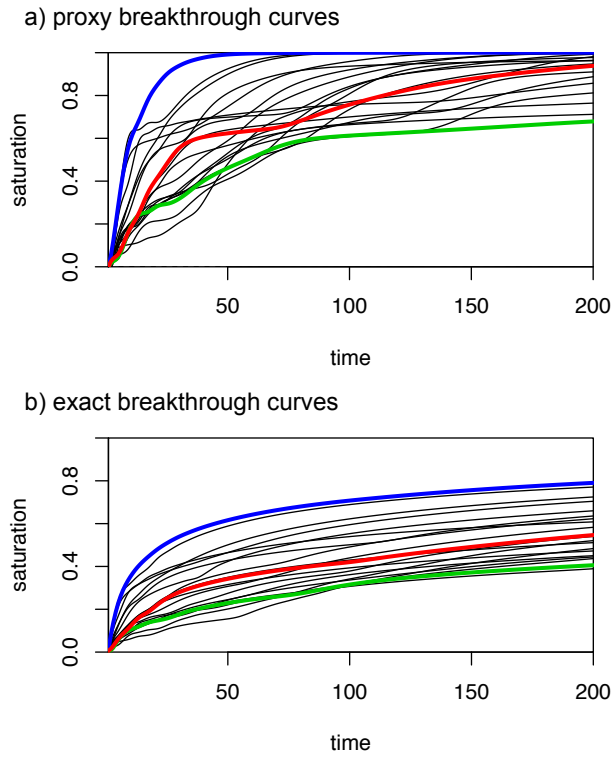


Figure 4. The learning set: (a) proxy curves and (b) exact curves recast as functional objects for the $N_l = 20$ realizations in the learning set. The thicker blue curves correspond to the realization in Fig. 3a), the red curves to 3b), and the green curves to 3c).

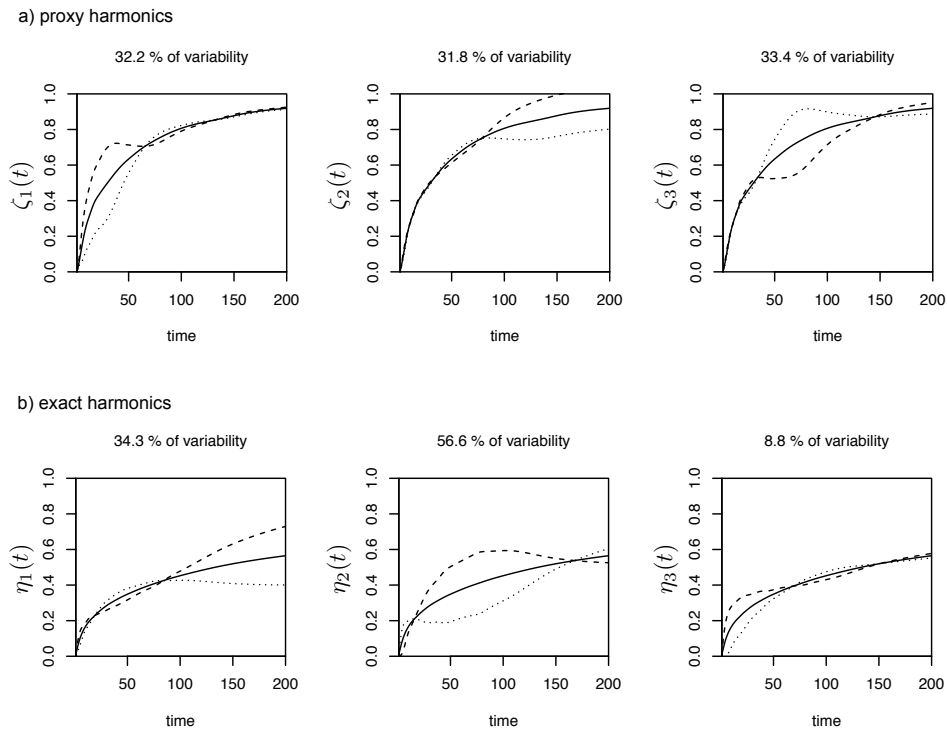


Figure 5. The three first rotated functional principal components (harmonics) extracted from the learning set are plotted for the proxy curves (top) and for the exact curves (bottom). The solid line is the mean curve and the dotted lines represent the variability around the mean described by the corresponding harmonic.

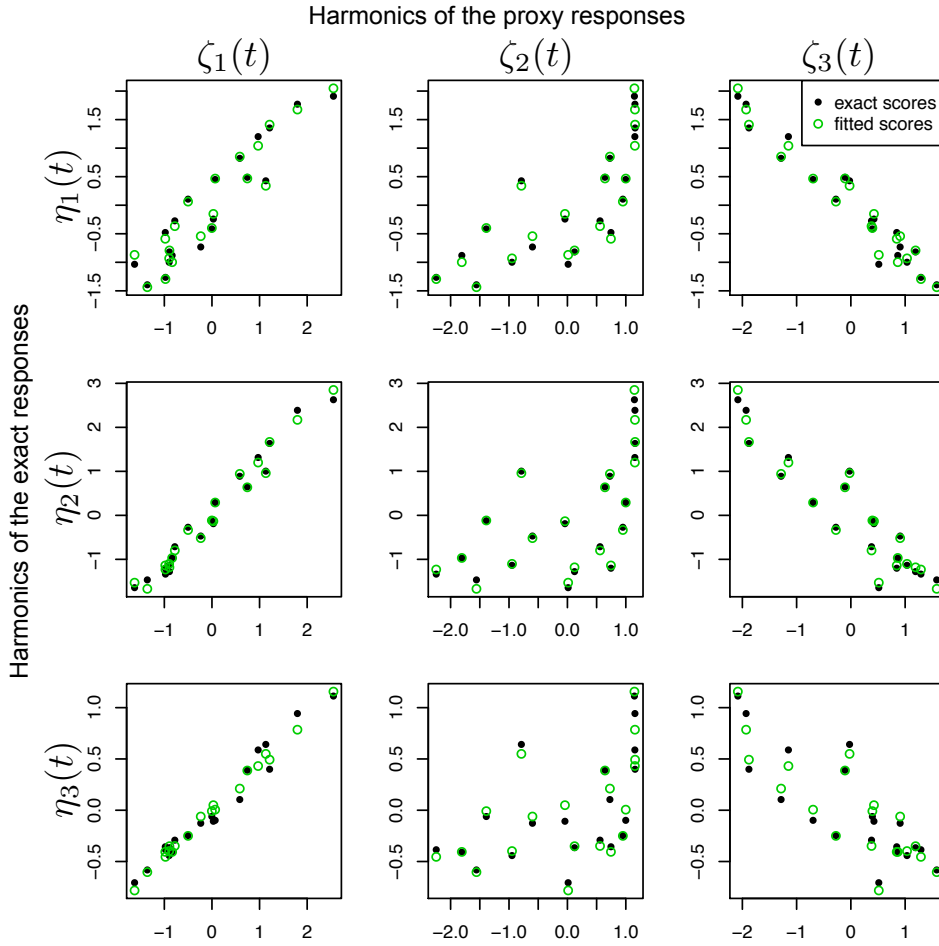


Figure 6. The scores, with respect to the first three harmonics $\{\eta_i(t)\}_{i=1,2,3}$, of the exact curves are plotted as functions of the scores for the approximate curves with respect to the harmonics $\{\zeta_i(t)\}_{i=1,2,3}$. The filled (black) circles correspond to the exact score, the empty circles (green) to the prediction of the scores by the OSL linear regression. The visualization is helpful to assess whether the linear regression model describes the relationship between proxy and exact curves in the learning set.

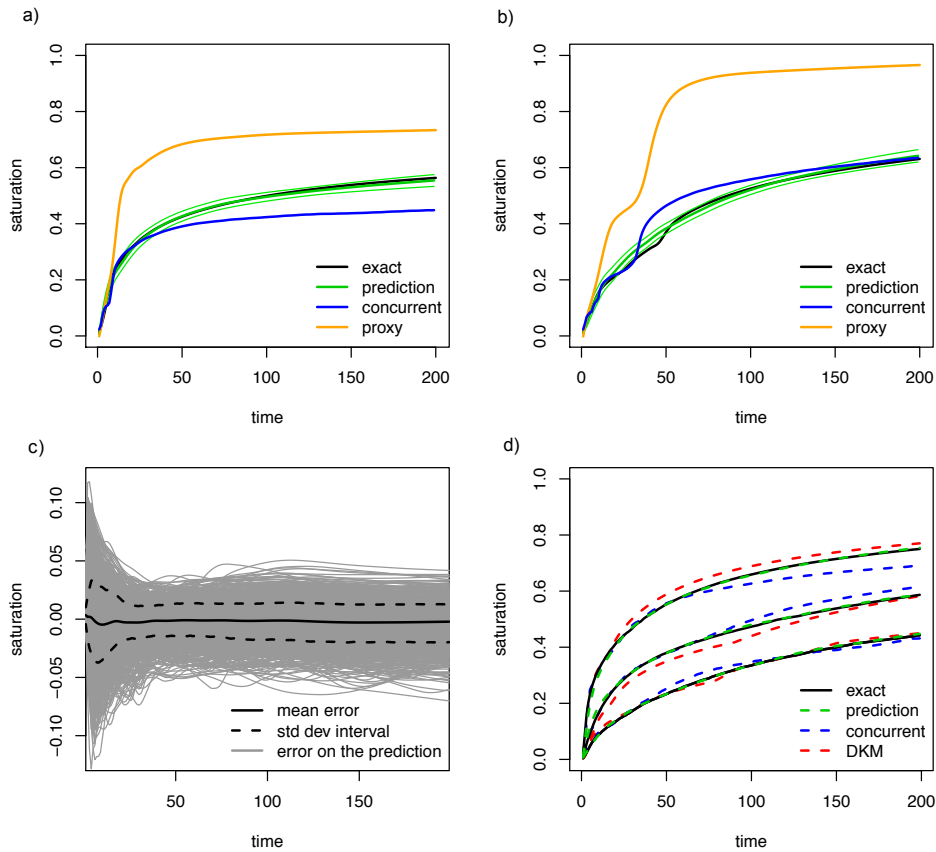
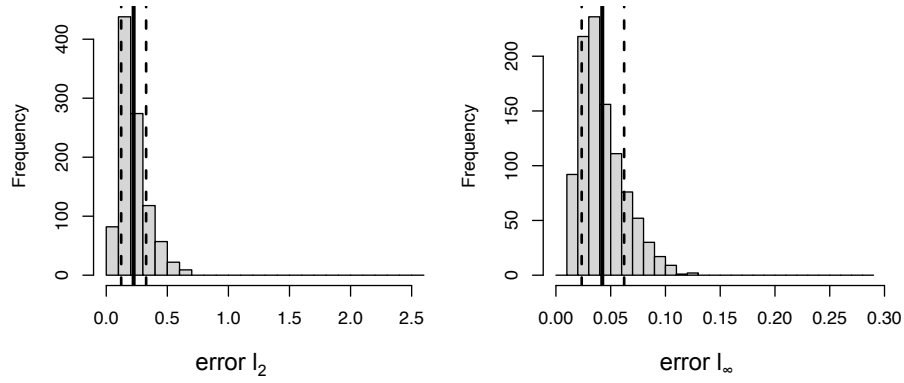


Figure 7. (a and b) the predicted responses (with 2σ -confidence intervals) of two realizations that are not in the learning set. (c) Prediction error of all $N_r = 1000$ realizations (gray curves), the mean error (continuous line), and the mean \pm one standard deviation (dotted lines) are represented. (d) P10, P50 and P90 quantiles curves obtained with the different models and compared to the reference quantile curves computed using the whole set of exact responses (solid black line).

a) FPCA error model



b) Concurrent model

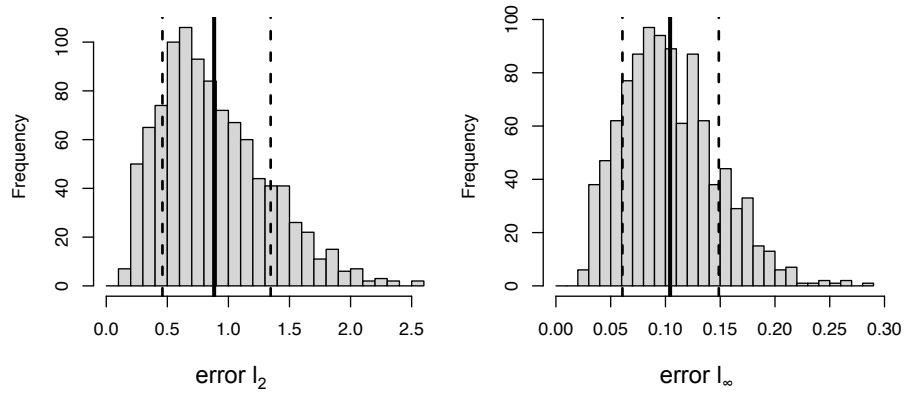


Figure 8. Histograms of the distribution of the l_2 error (left) and l_∞ error (right), (a) for the predictions of the FPCA model and (b) for the predictions of the concurrent model. The mean (continuous line) together with the mean \pm one standard deviation (dotted lines) are represented.

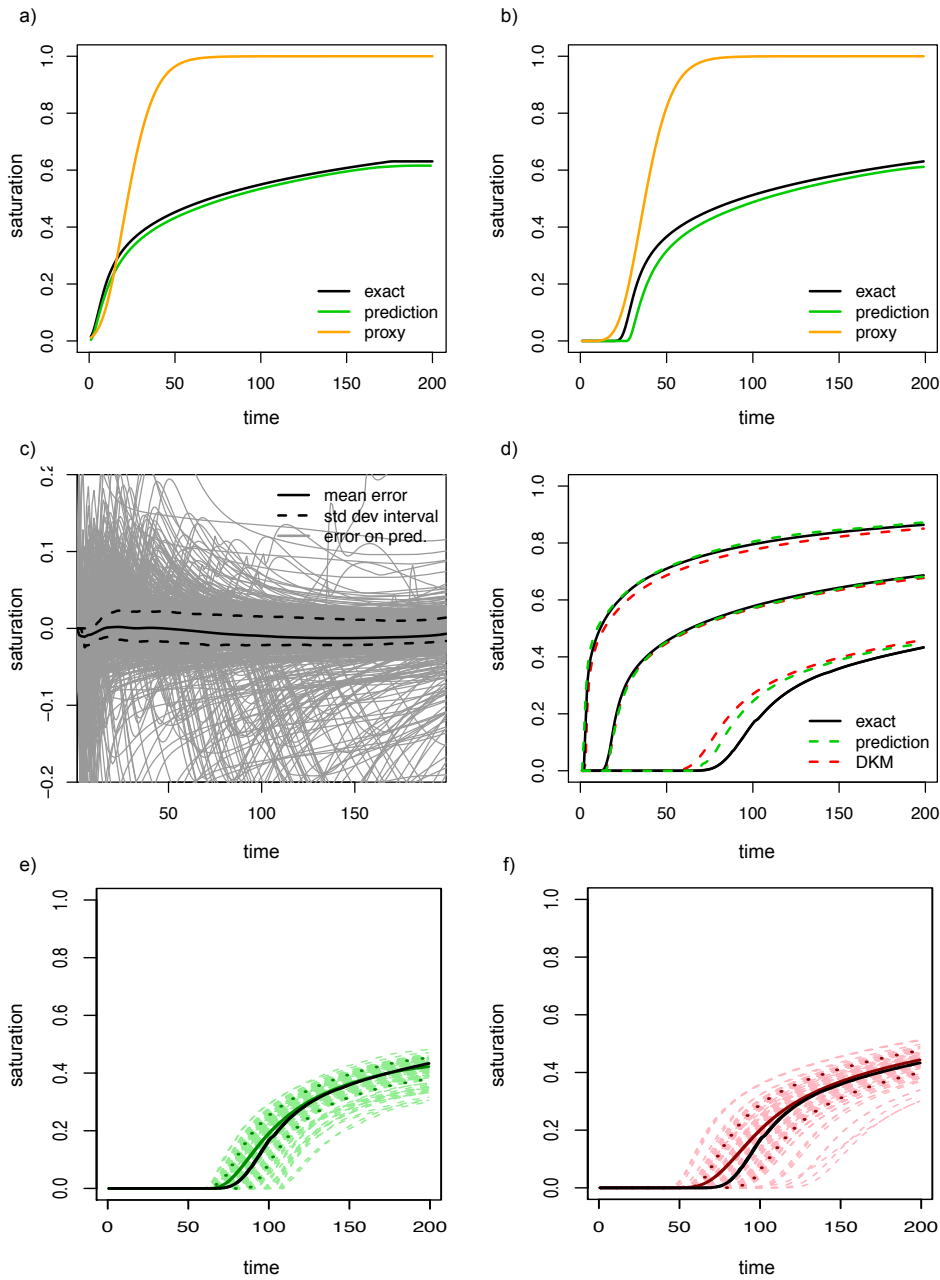


Figure 9. (a and b) predicted responses (before and after translation) of a realization that is not in the learning set. (c) Prediction error of all $N_r = 1000$ realizations (grey curves), the mean error (continuous line), and the mean \pm one standard deviation (dotted lines) are represented. (d) P10, P50 and P90 quantiles curves obtained with the different models and compared to the reference quantile curves computed using the whole set of exact responses (solid black line). (e), respectively (f), shows the P10 FPCA, respectively DKM, predictions of the P10 quantile for the 200 clusterings.

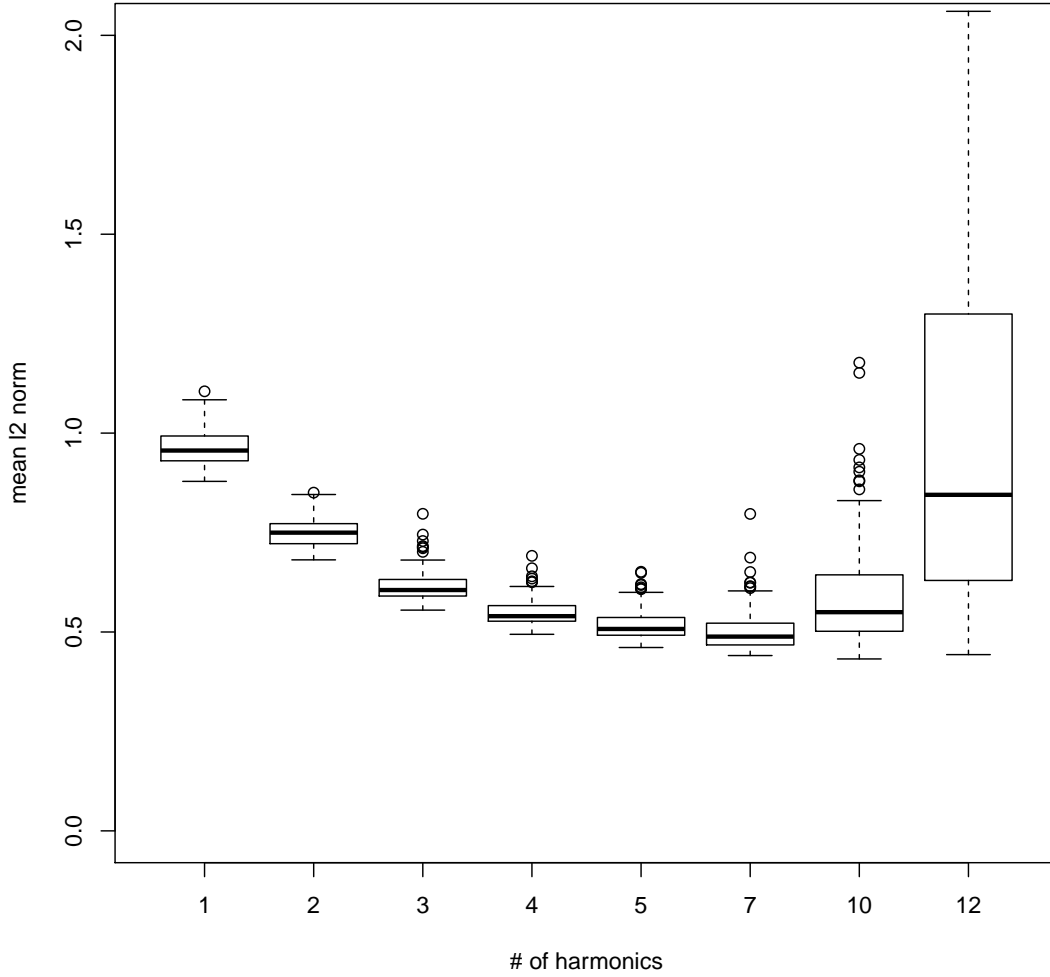


Figure 10. Boxplots of the prediction error (calculated as mean l_2 norm of the error of the predicted curves) as a function of the number of harmonics used to describe the proxy and exact curves in the learning set. The boxplots represent the statistics of the prediction errors over 200 clusterings in function of the number of harmonics. The thick line indicates the median error; the box the 1σ interval; the bars the 2σ interval; and the circles are the outliers (for 12 harmonics they are out of scale).