



HAL
open science

Recursive co-kriging model for Design of Computer experiments with multiple levels of fidelity with an application to hydrodynamic

Loic Le Gratiet

► **To cite this version:**

Loic Le Gratiet. Recursive co-kriging model for Design of Computer experiments with multiple levels of fidelity with an application to hydrodynamic. 2012. hal-00737332v2

HAL Id: hal-00737332

<https://hal.science/hal-00737332v2>

Preprint submitted on 29 Oct 2012 (v2), last revised 10 Jan 2013 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Recursive co-kriging model for Design of Computer experiments with multiple levels of fidelity with an application to hydrodynamic

Loic Le Gratiet
CEA, DAM, DIF, F-91297 Arpajon, France
loic.le-gratiet@cea.fr

October 29, 2012

In many practical cases, a sensitivity analysis or an optimization of a complex time consuming computer code requires to build a fast running approximation of it - also called surrogate model. We consider in this paper the problem of building a surrogate model of a complex computer code which can be run at different levels of accuracy. The co-kriging based surrogate model is a promising tool to build such an approximation. The idea is to improve the surrogate model by using fast and less accurate versions of the code. We present here a new approach to perform a multi-fidelity co-kriging model which is based on a recursive formulation. The strength of this new method is that the co-kriging model is built through a series of independent kriging models. From them, some properties of classical kriging models can naturally be extended to the presented co-kriging model such as a fast cross-validation procedure. Moreover, based on a Bayes linear formulation, an extension of the universal kriging equations are provided for the co-kriging model. Finally, the proposed model has the advantage to reduce the computational complexity compared to the previous models. The multi-fidelity model is successfully applied to emulate a hydrodynamic simulator. This real example illustrates the efficiency of the recursive model.

Keywords. surrogate models, universal co-kriging, recursive model, fast cross-validation, multi-fidelity computer code.

1 Introduction

Computer codes are widely used in science to describe physical phenomena. Advances in physics and computer science lead to increased complexity for the simulators. Therefore, it is common for the physicist to have different versions of a code which have different levels of accuracy and cost. Usually, to design and analyze a complex computer code, a fast approximation of it - also call surrogate model - is built in order to avoid prohibitive computational cost.

A very popular method to build surrogate model is the Gaussian process regression, also named kriging model. It is a particular class of surrogate models which makes the assumption that the response of the complex code is a realization of a Gaussian process. This method was originally introduced in used in geostatistics by [Krige, 1951] and [Matheron, 1963] and it was then proposed in the field of computer experiments by [Sacks et al., 1989]. During the

last decades, this method has become widely used and investigated. The reader is referred to the books of [Stein, 1999], [Santner et al., 2003] and [Rasmussen and Williams, 2006] for more detail about it.

A question of interest is how to build a predictive model using data from experiments of multiple levels of fidelity. Indeed, complex computer codes can be extremely expensive and sometimes we cannot have, under reasonable time constraint, enough simulations to sample the input parameter space with enough density and range. In this case, it could be worth using fast versions of the code (which can be old or coarse versions of it) to improve its approximation. Our objective is hence to build a multi-fidelity surrogate model which is able to use the information obtained from the fast versions of the code. Such models have been presented in the literature [Craig et al., 1998], [Kennedy and O’Hagan, 2000], [Forrester et al., 2007], [Qian and Wu, 2008] and [Cumming and Goldstein, 2009].

The first multi-fidelity model proposed in [Craig et al., 1998] is based on a linear regression formulation. Then [Cumming and Goldstein, 2009] have improved this model by using a Bayes linear formulation. The reader is referred to [Goldstein and Wooff, 2007] for further details about the Bayes linear approach. The methods suggested by [Craig et al., 1998] and [Cumming and Goldstein, 2009] have the strength to be relatively not computationally expensive but as it is based on a linear regression formulation, it could suffer from a lack of accuracy. Another approach is to use an extension of kriging for multiple response models which is called co-kriging. The idea was implemented by [Kennedy and O’Hagan, 2000] who present a co-kriging model based on an autoregressive relation between the different code levels. This method has become very popular and many authors have developed it. In particular, [Forrester et al., 2007] presents the use of co-kriging for multi-fidelity optimization and [Qian and Wu, 2008] proposed a Bayesian formulation of it.

The strength of the co-kriging model is that it gives very good predictive models but it is often computationally expensive, especially when the number of simulations is large. Furthermore, large data set can generate problems such as ill-conditioned covariance matrices. It is even more difficult to deal with these problems for co-kriging since the total number of observations is the sum of the observations at all code levels.

In this paper, we adopt a new approach for multi-fidelity surrogate modeling which uses a co-kriging model but with a recursive formulation. In fact, our model is able to build a s -level co-kriging model by building s independent krigings. This approach significantly reduces the complexity of the model since it divides the total number of observations on groups of observations corresponding to the ones of each level. Therefore, we will have s sub-matrices to invert which is less expensive than a big one and the estimation of the parameters can be performed separately. Finally, one of the main strengths of this approach is that it allows us to naturally extend classical results of kriging to the considered co-kriging model. In particular, we generalize and adapt the equations of the fast cross-validation proposed by [Dubrule, 1983] and we propose an universal co-kriging which is the natural extension of the well known universal kriging equations [Matheron, 1969].

2 Multi-fidelity Gaussian process regression.

In a first subsection, we briefly present a first approach to build multi-fidelity model suggested by [Kennedy and O’Hagan, 2000] that uses a co-kriging model. In the next subsection, we detail our recursive approach to build a multi-fidelity recursive model. The recursive formulation

of the multi-fidelity model is the first novelty of this paper. We will see in the next sections that the new formulation allows us to find very important results about the co-kriging model and to reduce its computational complexity. Furthermore, we prove that the two models are equivalent.

2.1 The classical autoregressive model.

Let us suppose that we have s levels of code $(z_t(x))_{t=1,\dots,s}$ sorted by increasing order of fidelity and modeled by Gaussian processes $(Z_t(x))_{t=1,\dots,s}$, $x \in Q$. We hence consider that $z_s(x)$ is the most accurate and costly code that we want to surrogate and $(z_t(x))_{t=1,\dots,s-1}$ are cheaper versions of it with $z_1(x)$ the less accurate one. We consider the following autoregressive model with $t = 2, \dots, s$:

$$\begin{cases} Z_t(x) = \rho_{t-1}(x)Z_{t-1}(x) + \delta_t(x) \\ Z_{t-1}(x) \perp \delta_t(x) \\ \rho_{t-1}(x) = g_{t-1}^T(x)\beta_{\rho_{t-1}} \end{cases} \quad (1)$$

where:

$$\delta_t(x) \sim \mathcal{GP}(f_t^T(x)\beta_t, \sigma_t^2 r_t(x, x')) \quad (2)$$

and:

$$Z_1(x) \sim \mathcal{GP}(f_1^T(x)\beta_1, \sigma_1^2 r_1(x, x')) \quad (3)$$

Here, T stands for the transpose, \perp denotes the orthogonality relationship, \mathcal{GP} designs a Gaussian Process, $g_{t-1}^T(x)$ is a vector of q_{t-1} regression functions, $f_t^T(x)$ is a vector of p_t regression functions, $r_t(x, x')$ is a correlation function, β_t is a p_t -dimensional vector, $\beta_{\rho_{t-1}}$ is a q_{t-1} -dimensional vector and σ_t^2 is a real. Since we suppose that the responses are realizations of Gaussian processes, the multi-fidelity model can be built by conditioning by the known responses of the codes at the different levels.

The previous model comes from the article of [Kennedy and O'Hagan, 2000]. It is induced by the following assumption: $\forall x \in Q$, if we know $Z_{t-1}(x)$, nothing more can be learned about $Z_t(x)$ from $Z_{t-1}(x')$ for $x \neq x'$.

Let us consider $\mathcal{Z}^{(s)} = (\mathcal{Z}_1^T, \dots, \mathcal{Z}_s^T)^T$ the Gaussian vector containing the values of the random processes $(Z_t(x))_{t=1,\dots,s}$ at the points in the experimental design sets $(D_t)_{t=1,\dots,s}$ with $D_s \subseteq D_{s-1} \subseteq \dots \subseteq D_1$ and $z^{(s)} = (z_1^T, \dots, z_s^T)^T$ a vector containing the values of $(z_t(x))_{t=1,\dots,s}$ at the points in $(D_t)_{t=1,\dots,s}$. The nested property of the experimental design sets is not necessary to build the model but it allows for a simple estimation of the model parameters. Since the codes are sorted in increasing order of fidelity it is not an unreasonable constraint for practical applications. By denoting $\beta = (\beta_1^T, \dots, \beta_s^T)^T$ the trend parameters, $\beta_\rho = (\beta_{\rho_1}^T, \dots, \beta_{\rho_{s-1}}^T)^T$ the trend of the adjustment parameters and $\sigma^2 = (\sigma_1^2, \dots, \sigma_s^2)$ the variance parameters, we have:

$$\forall x \in Q \quad [Z_s(x) | \mathcal{Z}^{(s)} = z^{(s)}, \beta, \beta_\rho, \sigma^2] \sim \mathcal{N}(m_{Z_s}(x), s_{Z_s}^2(x))$$

where:

$$m_{Z_s}(x) = h'_s(x)^T \beta + t_s(x)^T V_s^{-1} (z^{(s)} - H_s \beta) \quad (4)$$

and:

$$s_{Z_s}^2(x) = v_{Z_s}^2(x) - t_s(x)^T V_s^{-1} t_s(x) \quad (5)$$

The Gaussian process regression mean $m_{Z_s}(x)$ is the predictive model of the highest fidelity response $z_s(x)$ which is built with the known responses of all code levels $z^{(s)}$. The variance $s_{Z_s}^2(x)$ represents the predictive mean squared error of the model.

The matrix V_s is the covariance matrix of the Gaussian vector $\mathcal{Z}^{(s)}$, the vector $t_s(x)$ is the vector of covariance between $Z_s(x)$ and $\mathcal{Z}^{(s)}$, $H_s\beta$ is the mean of $\mathcal{Z}^{(s)}$, $h'_s(x)^T\beta$ is the mean of $Z_s(x)$ and $v_{Z_s}^2(x)$ is the variance of $Z_s(x)$. All these terms are built in terms of the experience vector at level t (6) and to the covariance between $Z_t(x)$ and $Z_{t'}(x')$ (7) and (8).

$$h'_t(x)^T = \left(\left(\prod_{i=1}^{t-1} \rho_i(x) \right) f_1^T(x), \left(\prod_{i=2}^{t-1} \rho_i(x) \right) f_2^T(x), \dots, \rho_{t-1}(x) f_{t-1}^T(x), f_t^T(x) \right) \quad (6)$$

Let us consider $t > t'$:

$$\text{cov}(Z_t(x), Z_{t'}(x') | \sigma^2, \beta, \beta_\rho) = \left(\prod_{i=t'}^{t-1} \rho_i(x) \right) \text{cov}(Z_{t'}(x), Z_{t'}(x') | \sigma^2, \beta, \beta_\rho) \quad (7)$$

with :

$$\text{cov}(Z_t(x), Z_t(x') | \sigma^2, \beta, \beta_\rho) = \sum_{j=1}^t \sigma_j^2 \left(\prod_{i=j}^{t-1} \rho_i(x) \rho_i(x') \right) r_j(x, x') \quad (8)$$

Remark. The model (1) is an extension of the model of [Kennedy and O'Hagan, 2000] in which the adjustment parameters $\rho_t(x)_{t=2, \dots, s}$ do not depend on x . We show in a practical application that this extension is worthwhile.

2.2 Recursive multi-fidelity model.

In this section, we present the new multi-fidelity model which is based on a recursive formulation. Let us consider the following model for $t = 2, \dots, s$:

$$\begin{cases} Z_t(x) = \rho_{t-1}(x) \tilde{Z}_{t-1}(x) + \delta_t(x) \\ \tilde{Z}_{t-1}(x) \perp \delta_t(x) \\ \rho_{t-1}(x) = g_{t-1}^T(x) \beta_{\rho_{t-1}} \end{cases} \quad (9)$$

where $\tilde{Z}_{t-1}(x)$ is a Gaussian process with distribution $[Z_{t-1}(x) | \mathcal{Z}^{(t-1)} = z^{(t-1)}, \beta_{t-1}, \beta_{\rho_{t-2}}, \sigma_{t-1}^2]$ and $D_s \subseteq D_{s-1} \subseteq \dots \subseteq D_1$. The unique difference with the previous model is that we express $Z_t(x)$ (the Gaussian process modeling the response at level t) as a function of the Gaussian process $Z_{t-1}(x)$ conditioned by the values $z^{(t-1)} = (z_1, \dots, z_{t-1})$ at points in the experimental design sets $(D_i)_{i=1, \dots, t-1}$. We note that, as in the previous model, the nested property is assumed to allow efficient estimations for the model parameters. The Gaussian processes $(\delta_t(x))_{t=2, \dots, s}$ have the same definition as previously and we have for $t = 2, \dots, s$:

$$[Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \beta_t, \beta_{\rho_{t-1}}, \sigma_t^2] \sim \mathcal{N}(\mu_{Z_t}(x), s_{Z_t}^2(x)) \quad (10)$$

where:

$$\mu_{Z_t}(x) = \rho_{t-1}(x) \mu_{Z_{t-1}}(x) + f_t^T(x) \beta_t + r_t^T(x) R_t^{-1} (z_t - \rho_{t-1}(D_t) \odot z_{t-1}(D_t) - F_t \beta_t) \quad (11)$$

and:

$$\sigma_{Z_t}^2(x) = \rho_{t-1}^2(x) \sigma_{Z_{t-1}}^2(x) + \sigma_t^2 (1 - r_t^T(x) R_t^{-1} r_t(x)) \quad (12)$$

The notation \odot represents the element by element matrix product. R_t is the correlation matrix $R_t = (r_t(x, x'))_{x, x' \in D_t}$ and $r_t^T(x)$ is the correlation vector $r_t^T(x) = (r_t(x, x'))_{x' \in D_t}$. We denote by $\rho_t(D_{t-1})$ the vector containing the values of $\rho_t(x)$ for $x \in D_{t-1}$, $z_t(D_{t-1})$ the vector containing the known values of $Z_t(x)$ at points in D_{t-1} and F_t is the experience matrix containing the values of $f_t(x)^T$ on D_t .

The mean $\mu_{Z_t}(x)$ is the surrogate model of the response at level t , $1 \leq t \leq s$, taking into account the known values of the t first levels of responses $(z_i)_{i=1, \dots, t}$ and the variance $\sigma_{Z_t}^2(x)$ represents the mean squared error of this model. The mean and the variance of the Gaussian process regression at level t being expressed in function of the ones of level $t-1$, we so have a recursive multi-fidelity metamodel. Furthermore, in this new formulation, it is clearly emphasized that the mean of the predictive distribution does not depend on the variance parameters $(\sigma_t^2)_{t=1, \dots, s}$. This is a classical result of kriging model which states that for covariance kernels of the form $k(x, x') = \sigma^2 r(x, x')$, the mean of the kriging model is independent of σ^2 .

Remark. The previous comment highlights an important strength of the recursive formulation. Indeed, contrary to the formulation suggested in [Kennedy and O'Hagan, 2000], once the multi-fidelity model is built, it provides the surrogate models of all the responses $(z_t(x))_{t=1, \dots, s}$.

Furthermore, from this formulation, we can directly deduce that building a s -level co-kriging is equivalent to build s independent krigings. This implies a reduction of the model complexity. Indeed, the inversion of the matrix V_s of size $\sum_{i=1}^s n_i \times \sum_{i=1}^s n_i$ is more expensive than the inversions of s matrices $(R_t)_{t=1, \dots, s}$ of size $(n_t \times n_t)_{t=1, \dots, s}$ where n_t corresponds to the size of the vector z_t at level $t = 1, \dots, s$. We also reduce the memory cost since storing the matrix V_s required more memory than storing the s matrices $(R_t)_{t=1, \dots, s}$. Then, we note that the model with this formulation is more interpretable since we can deduce the impact of each level of response into the model error through $(\sigma_{Z_t}^2(x))_{t=1, \dots, s}$. Finally we will see that it allows us to adapt classical kriging results to the multi-fidelity co-kriging model (e.g. universal kriging and fast cross-validation).

We have the following proposition.

Proposition 1 *Let us consider s Gaussian processes $(Z_t(x))_{t=1, \dots, s}$ and $\mathcal{Z}^{(s)} = (Z_t)_{t=1, \dots, s}$ the Gaussian vector containing the values of $(Z_t(x))_{t=1, \dots, s}$ at points in $(D_t)_{t=1, \dots, s}$ with $D_s \subseteq D_{s-1} \subseteq \dots \subseteq D_1$. If we consider the mean and the variance (4) and (5) induced by the model (1) when we condition the Gaussian process $Z_s(x)$ by the known values $z^{(s)}$ of $\mathcal{Z}^{(s)}$ and the mean and the variance (11) and (12) induced by the model (16) when we condition $Z_s(x)$ by $z^{(s)}$, then, we have:*

$$\begin{aligned} \mu_{Z_s}(x) &= m_{Z_s}(x) \\ \sigma_{Z_s}^2(x) &= s_{Z_s}^2(x) \end{aligned}$$

The proof of the proposition is given in Appendix A.1. It shows that the model of [Kennedy and O'Hagan, 2000] and the recursive model (16) have the same mean and covariance function. Therefore, predictive distributions of the two models are identical and the recursive model has the same strengths as the one of [Kennedy and O'Hagan, 2000] to which we add the benefits mentioned in the previous remark.

2.3 Parameter estimations

We present in this section a Bayesian estimation of the parameter $\psi = (\beta, \beta_\rho, \sigma^2)$ focusing on conjugate and non-informative distributions for the priors. This allows us to obtain closed form expressions for the estimations of the parameters. Furthermore, from the non-informative case, we can obtain the estimates given by a maximum likelihood method. The presented formulas can hence be used in a frequentist approach. We note that the recursive formulation directly shows us that the estimations of the parameters $(\beta_t, \beta_{\rho_{t-1}}, \sigma_t^2)_{t=1, \dots, s}$ and (β_1, σ_1^2) can be performed separately.

We use in this section the notation **info** to design the case where all the priors are informative and **ninfo** to design the case where all the priors are non-informative. It would be possible to address the case of a mixture of informative and non-informative priors. For the non-informative case, we use the "Jeffreys priors" [Jeffreys, 1961]:

$$p(\beta_1 | \sigma_1^2) \propto 1, \quad p(\sigma_1^2) \propto \frac{1}{\sigma_1^2}, \quad p(\beta_{\rho_{t-1}}, \beta_t | z^{(t-1)}, \sigma_t^2) \propto 1, \quad p(\sigma_t^2 | z^{(t-1)}) \propto \frac{1}{\sigma_t^2} \quad (13)$$

where $t = 2, \dots, s$. For the informative case, we consider the following conjugate prior distributions:

$$\begin{aligned} [\beta_1 | \sigma_1^2] &\sim \mathcal{N}_{p_1}(b_1, \sigma_1^2 V_1) \\ [\beta_{\rho_{t-1}}, \beta_t | z^{(t-1)}, \sigma_t^2] &\sim \mathcal{N}_{q_{t-1} + p_t} \left(b_t = \begin{pmatrix} b_{t-1}^\rho \\ b_t^\beta \end{pmatrix}, \sigma_t^2 V_t = \sigma_t^2 \begin{pmatrix} V_{t-1}^\rho & 0 \\ 0 & V_t^\beta \end{pmatrix} \right) \\ [\sigma_1^2] &\sim \mathcal{IG}(\alpha_1, \gamma_1), \quad [\sigma_t^2 | z^{(t-1)}] \sim \mathcal{IG}(\alpha_t, \gamma_t) \end{aligned}$$

with b_1 a vector a size p_1 , b_{t-1}^ρ a vector of size q_{t-1} , b_t^β a vector of size p_t , V_1 a $p_1 \times p_1$ matrix, V_{t-1}^ρ a $q_{t-1} \times q_{t-1}$ matrix, V_t^β a $p_t \times p_t$ matrix and $\alpha_1, \gamma_1, \alpha_t, \gamma_t > 0$. These informative priors allow the user to prescribe the means and the variances of all parameters. The choice of conjugate priors allows us to have closed form expressions for the parameter estimations. Indeed, we have:

$$[\beta_1 | z_1, \sigma_1^2] \sim \mathcal{N}_{p_1}(\Sigma_1 \nu_1, \Sigma_1) \quad [\beta_{\rho_{t-1}}, \beta_t | z^{(t)}, \sigma_t^2] \sim \mathcal{N}_{q_{t-1} + q_t}(\Sigma_t \nu_t, \Sigma_t) \quad (14)$$

where, for $t \geq 1$:

$$\Sigma_t = \begin{cases} [H_t^T \frac{R_t^{-1}}{\sigma_2^2} H_t + \frac{V_t^{-1}}{\sigma_2^2}]^{-1} & \mathbf{info} \\ [H_t^T \frac{R_t^{-1}}{\sigma_2^2} H_t]^{-1} & \mathbf{ninfo} \end{cases} \quad \nu_t = \begin{cases} [H_t^T \frac{R_t^{-1}}{\sigma_2^2} z_t + \frac{V_t^{-1}}{\sigma_2^2} b_t] & \mathbf{info} \\ [H_t^T \frac{R_t^{-1}}{\sigma_2^2} z_t] & \mathbf{ninfo} \end{cases} \quad (15)$$

with $H_1 = F_1$ and for $t > 1$, $H_t = [G_{t-1} \odot (z_{t-1}(D_t) \mathbf{1}_{q_{t-1}}^T) \quad F_t]$ where G_{t-1} is the experience matrix containing the values of $g_{t-1}(x)^T$ in D_t . Furthermore, we have for $t \geq 1$:

$$[\sigma_t^2 | z^{(t)}] \sim \mathcal{IG}(a_t, \frac{Q_t}{2}) \quad (16)$$

where:

$$Q_t = \begin{cases} \gamma_t + (b_t - \hat{\lambda}_t)^T (V_t + [H_t^T R_t^{-1} H_t]^{-1})^{-1} (b_t - \hat{\lambda}_t) + \hat{Q}_t & \mathbf{info} \\ \hat{Q}_t & \mathbf{ninfo} \end{cases}$$

with $\hat{Q}_t = (z_t - H_t \hat{\lambda}_t)^T R_t^{-1} (z_t - H_t \hat{\lambda}_t)$, $\hat{\lambda}_t = (H_t^T R_t^{-1} H_t F)^{-1} H_t^T R_t^{-1} z_t$ and :

$$a_t = \begin{cases} \frac{n_t}{2} + \alpha_t & \mathbf{info} \\ \frac{n_t - p_t - q_{t-1}}{2} & \mathbf{ninfo} \end{cases}$$

with the convention $q_0 = 0$.

We highlight that the maximum likelihood estimates for the parameters β_1 and $(\beta_{\rho_{t-1}}, \beta_t)$ are given by the means of the posterior distributions of the Bayesian estimations in the non-informative case. Furthermore, the restricted maximum likelihood estimate of the variance parameter σ_t^2 can also be deduced from the posterior distribution of the Bayesian estimation in the non-informative case and is given by $\hat{\sigma}_{t,\text{EML}}^2 = \frac{Q_t}{2a_t}$. The restricted maximum likelihood estimation is a method which allows us to reduce the bias of the maximum likelihood estimation [Patterson and Thompson, 1971].

3 Universal co-kriging model

We see in equation (10) that the predictive distribution of $Z_s(x)$ is conditioned by the observations z and the parameters β , β_ρ and σ^2 . The objective of a Bayesian prediction is to integrate the uncertainty due to the parameter estimations into the predictive distribution. Indeed, in the previous subsection, we have expressed the posterior distributions of the variance parameters $(\sigma_t^2)_{t=1,\dots,s}$ conditionally to the observations and the posterior distributions of the trend parameters β_1 and $(\beta_{\rho_{t-1}}, \beta_t)_{t=2,\dots,s}$ conditionally to the observations and the variance parameters. Thus, using the Bayes formula, we can easily obtain a predictive distribution only conditioned by the observations by integrating into it the posterior distributions of the parameters.

Nonetheless, this predictive distribution is clearly not Gaussian and can be expensive to obtain. In particular, we cannot have a closed form expression for the predictive distribution. Therefore, it is relevant to consider in our analysis only the mean $\mathbb{E}[Z_s(x)|\mathcal{Z}^{(s)} = z^{(s)}]$ and the variance $\text{Var}(Z_s(x)|\mathcal{Z}^{(s)} = z^{(s)})$.

The following proposition giving the closed form expressions of the mean and the variance of the predictive distribution only conditioned by the observations is a novelty. The proof of this proposition is based on the recursive formulation which emphasizes the strength of this new approach.

Proposition 2 *Let us consider s Gaussian processes $(Z_t(x))_{t=1,\dots,s}$ and $\mathcal{Z}^{(s)} = (Z_t)_{t=1,\dots,s}$ the Gaussian vector containing the values of $(Z_t(x))_{t=1,\dots,s}$ at points in $(D_t)_{t=1,\dots,s}$ with $D_s \subseteq D_{s-1} \subseteq \dots \subseteq D_1$. If we consider the conditional predictive distribution in equation (10) and the posterior distributions of the parameters given in equations (14) and (16), then we have for $t = 1, \dots, s$:*

$$\mathbb{E}[Z_t(x)|\mathcal{Z}^{(t)} = z^{(t)}] = h_t^T(x) \Sigma_t \nu_t + r_t^T(x) R_t^{-1} (z_t - H_t \Sigma_t \nu_t) \quad (17)$$

with $h_1^T = f_1^T$, $H_1 = F_1$ and for $t > 1$, $h_t^T(x) = (g_{t-1}(x)^T \mathbb{E}[Z_{t-1}(x)|\mathcal{Z}_{t-1} = z_{t-1}] \quad f_t^T(x))$ and $H_t = [G_{t-1} \odot (z_{t-1}(D_t) \mathbf{1}_{q_{t-1}}^T) \quad F_t]$. Furthermore, we have:

$$\begin{aligned} \text{Var}(Z_t(x)|\mathcal{Z}^{(t)} = z^{(t)}) &= \hat{\rho}_{t-1}^2(x) \text{Var}(Z_{t-1}(x)|\mathcal{Z}^{(t-1)} = z^{(t-1)}) + \frac{Q_t}{2(a_t-1)} (1 - r_t^T(x) R_t^{-1} r_t^T(x)) \\ &\quad + (h_t^T - r_t^T(x) R_t^{-1} H_t) \Sigma_t (h_t^T - r_t^T(x) R_t^{-1} H_t)^T \end{aligned} \quad (18)$$

with $\hat{\rho}_{t-1}(x) = [\Sigma_t \nu_t]_{1,\dots,q_{t-1}}$.

The proof of Proposition 2 is given in Appendix A.2. We note that, in the mean of the predictive distribution, the parameters have been replaced by their posterior means. Furthermore, in the variance of the predictive distribution, the variance parameter has been replaced by its posterior mean and the term $(h_t^T - r_t^T(x)R_t^{-1}H_t) \Sigma_t (h_t^T - r_t^T(x)R_t^{-1}H_t)^T$ has been added. It represents the uncertainty due to the estimation of the regression parameters (including the adjustment coefficient). We call these formulas the universal co-kriging equations due to their similarities with the well-known universal kriging equations (they are identical for $s = 1$).

4 Fast cross-validation for kriging and co-kriging surrogate models

The idea of a cross-validation procedure is to split the experimental design set into two disjoint sets, one is used for training and the other one is used to monitor the performance of the surrogate model. The idea is that, the performance on the test set can be used as a proxy for the generalization error. A particular case of this method is the Leave-One-Out Cross-Validation (noted LOO-CV) where n test sets are obtained by removing one observation at a time. This procedure can be time-consuming for a kriging model but [Dubrule, 1983], [Rasmussen and Williams, 2006] and [Zhang and Wang, 2009] show that there are computational shortcuts. We present in this section their adaptation for co-kriging models. Furthermore, the cross-validation equations proposed in this section extend the previous ones even for $s = 1$ (i.e. the classical kriging model) since they do not suppose that the regression and the variance coefficients are known. Therefore, those parameters are re-estimated for each training set. We note that the re-estimation of the variance coefficient is a novelty which is important since fixing this parameter can lead to big errors for the estimation of the cross-validation predictive variance when the number of observations is small or when the number of points in the test set is important.

If we denote by ξ_s the set of indices of n_{test} points in D_s constituting the test set D_{test} and ξ_t , $1 \leq t < s$, the corresponding set of indices in D_t - indeed, we have $D_s \subset D_{s-1} \subset \dots \subset D_1$, therefore $D_{test} \subset D_t$. The nested experimental design assumption implies that, in the cross-validation procedure, if we remove a group of points from D_s we can also remove it from D_t , $1 \leq t \leq s$.

The following proposition gives the vectors of the cross-validation predictive errors and variances at points in the test set D_{test} when we remove them from the t highest levels of code. In the proposition, we consider that we are in the non-informative case for the parameter estimation (see Section 2.3) but it can be easily extended to the informative case presented in Section 2.3.

Notations: If ξ is a set of indices, then $A_{[\xi, \xi]}$ is the sub-matrix of elements $\xi \times \xi$ of A , $a_{[\xi]}$ is the sub-vector of elements ξ of a , $B_{[-\xi]}$ represents the matrix B minus the rows of index ξ , $C_{[-\xi, -\xi]}$ is the sub-matrix of C in which we remove the elements of index $-\xi \times -\xi$ and $C_{[-\xi, \xi]}$ is the sub-matrix of C in which we remove the rows of index ξ and keep the columns of index ξ .

Proposition 3 *Let us consider s Gaussian processes $(Z_t(x))_{t=1, \dots, s}$ and $\mathcal{Z}^{(s)} = (Z_t)_{t=1, \dots, s}$ the Gaussian vector containing the values of $(Z_t(x))_{t=1, \dots, s}$ at points in $(D_t)_{t=1, \dots, s}$ with $D_s \subseteq D_{s-1} \subseteq \dots \subseteq D_1$. We note D_{test} a set made with the points of index ξ_s of D_s and ξ_t the*

corresponding points in D_t with $1 \leq t \leq s$. Then, if we note ε_{Z_s, ξ_s} the errors (i.e. real values minus predicted values) of the cross-validation procedure when we remove the points of D_{test} from the t highest levels of code, we have:

$$(\varepsilon_{Z_s, \xi_s} - \rho_{s-1}(D_{test}) \odot \varepsilon_{Z_{s-1}, \xi_{s-1}}) [R_s^{-1}]_{[\xi_s, \xi_s]} = [R_s^{-1} (z_s - H_s \lambda_{s, -\xi_s})]_{[\xi_s]} \quad (19)$$

with $\varepsilon_{Z_u, \xi_u} = 0$ when $u < t$, $\lambda_{s, -\xi_s} ([H_s^T]_{[-\xi_s]} K_s [H_s]_{[-\xi_s]}) = [H_s^T]_{[-\xi_s]} K_s z_s (D_s \setminus D_{test})$ and:

$$K_s = [R_s^{-1}]_{[-\xi_s, -\xi_s]} - [R_s^{-1}]_{[-\xi_s, \xi_s]} \left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1} [R_s^{-1}]_{[\xi_s, -\xi_s]} \quad (20)$$

Furthermore, if we note σ_{Z_s, ξ_s}^2 the variances of the corresponding cross-validation procedure, we have:

$$\sigma_{Z_s, \xi_s}^2 = \rho_{s-1}^2(D_{test}) \odot \sigma_{Z_{s-1}, \xi_{s-1}}^2 + \sigma_{s, -\xi_s}^2 \text{diag} \left(\left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1} \right) + \mathcal{V}_s \quad (21)$$

with:

$$\sigma_{s, -\xi_s}^2 = \frac{(z_s(D_s \setminus D_{test}) - [H_s]_{[-\xi_s]} \lambda_{s, -\xi_s})^T K_s (z_s(D_s \setminus D_{test}) - [H_s]_{[-\xi_s]} \lambda_{s, -\xi_s})}{n_s - p_s - q_{s-1} - n_{train}} \quad (22)$$

where $\sigma_{u, -\xi_u}^2 = 0$ when $u < t$, n_{train} is the length of the index vector ξ_s , $H_s = [G_{s-1} \odot (z_{s-1}(D_s) \mathbf{1}_{q_{s-1}}^T \quad F_s)]$ and:

$$\mathcal{V}_s = \mathcal{U}_s^T ([H_s^T]_{[-\xi_s]} K_s [H_s]_{[-\xi_s]})^{-1} \mathcal{U}_s \quad (23)$$

with $\mathcal{U}_s = \left(([R_s^{-1}]_{[\xi_s, \xi_s]})^{-1} [R_s^{-1} H_s]_{[\xi_s]} \right)$.

We note that these equations are also valid when $s = 1$, i.e. for kriging model. We hence have closed form expressions for the equations of a k -fold cross-validation with a re-estimation of the regression and variance parameters. These expressions can be deduced from the universal co-kriging equations. The complexity of this procedure is essentially determined by the inversion of the matrices $\left([R_u^{-1}]_{[\xi_u, \xi_u]} \right)_{u=t, \dots, s}$ of size $n_{test} \times n_{test}$. Furthermore, if we suppose the parameters of variance and/or trend as known, we do not have to compute $\sigma_{t, -\xi_t}^2$ and/or $\lambda_{t, -\xi_t}$ (they are fixed to their estimated value, i.e. $\sigma_{t, -\xi_t}^2 = \frac{Q_t}{2(a_t - 1)}$ and $\lambda_{t, -\xi_t} = \Sigma_t \nu_t$, see Section 2.3) which reduces substantially the complexity of the method. These equations generalize those of [Dubrule, 1983] and [Zhang and Wang, 2009] where the variance $\sigma_{t, -\xi_t}^2$ is supposed to be known. Finally, the term \mathcal{V}_s corresponds to the added term due to the parameter estimations in the universal co-kriging. Therefore, if the trend parameters are supposed known, this term is equal to 0. The proof of Proposition 3 is given in Appendix A.3.

5 Application: hydrodynamic simulator

In this section we apply our co-kriging method to the hydrodynamic code ‘‘MELTEM’’. This code simulates a second-order turbulence model for gaseous mixtures induced by Richtmyer-Meshkov instability [Grégoire et al., 2005]. Two input parameters x_1 and x_2 are considered. They are phenomenological coefficients used in the equations of the energy of dissipation

of the turbulent flow. These two coefficients vary in the region $[0.5, 1.5] \times [1.5, 2.3]$. The considered code outputs, called eps and L_c , are respectively the dissipation factor and the mixture characteristic length. The simulator is a finite-elements code which can be run at $s = 2$ levels of accuracy by altering the finite-elements mesh. The simple code $z_1(\cdot)$, using a coarse mesh, takes 15 seconds to produce an output whereas the complex code $z_2(\cdot)$, using a fine mesh, takes 8 minutes. The aim of the study is to build a prediction as accurate as possible using only a few runs of the complex code and to assess the uncertainty of this prediction. In particular, we use 5 runs for the complex code $z_2(x)$ and 25 runs for the cheap code $z_1(x)$. Then, we build an additional set of 175 points to test the accuracy of the models. Furthermore, no prior information is available: we are hence in the non-informative case.

5.1 Estimation of the hyper-parameters

In the previous sections, we have considered the correlation kernels $(r_t(x, x'))_{t=1, \dots, s}$ as known. In practical applications, we choose these kernels in a parameterized family of correlation kernels. Therefore, we consider kernels such that $r_t(x, x') = r_t(x, x'; \phi_t)$. The hyper-parameter ϕ_t can be estimated by maximizing the concentrated restricted log-likelihood [Santner et al., 2003] with respect to ϕ_t :

$$\log(|\det(R_t)|) + (n_t - p_t - q_{t-1}) \log(\sigma_{t,\text{EML}}^2) \quad (24)$$

with the convention $q_0 = 0$ and $\sigma_{t,\text{EML}}^2$ is the restricted likelihood estimate of the variance σ_t^2 (see Section 2.3). This minimization problem has to be solved numerically. It is a common choice to consider the hyper-parameters as known and to estimate them by maximum likelihood [Santner et al., 2003].

It is also possible to estimate the hyper-parameters $(\phi_t)_{t=1, \dots, s}$ by minimizing a loss function of a Leave-One-Out Cross-Validation procedure. Usually, the complexity of this procedure is $\mathcal{O}\left(\left(\sum_{i=1}^s n_i\right)^4\right)$. Nonetheless, thanks to Proposition 3, it is reduced to $\mathcal{O}\left(\sum_{i=1}^s n_i^3\right)$ since it is essentially determined by the inversions of the s matrices $(R_t^{-1})_{t=1, \dots, s}$. Therefore, the complexity for the estimation of $(\phi_t)_{t=1, \dots, s}$ is substantially reduced. Furthermore, the recursive formulation of the problem allows us to estimate the parameters $(\phi_t)_{t=1, \dots, s}$ one at a time by starting with ϕ_1 and estimating ϕ_t , $t = 2, \dots, s$, one after the other.

5.2 Comparison between kriging and multi-fidelity co-kriging

Before considering the real case study, we propose in this section a comparison between the kriging and co-kriging models when the number of runs n_2 for the complex code varies such that $n_2 = 5, 10, 15, 20, 25$. For the co-kriging model, we consider $n_1 = 25$ runs for the cheap code. In this section, we focus on the output eps .

To perform the comparison, we generate randomly 500 experimental design sets $(D_{2,i}, D_{1,i})_{i=1, \dots, 500}$ such that $D_{2,i} \subset D_{1,i}$, $i = 1, \dots, 500$, $D_{1,i}$ has n_1 points and $D_{2,i}$ has n_2 points.

We use for both kriging and co-kriging models a Matern $_{\frac{5}{2}}$ covariance kernel and we consider ρ , β_1 and β_2 as constant. The accuracies of the two models are evaluated on the test set composed of 175 observations. From them, the Root Mean Squared Error (RMSE) is computed: $\text{RMSE} = \left(\frac{1}{175} \sum_{i=1}^{175} (\mu_{Z_2}(x_i^{\text{test}}) - z_2(x_i^{\text{test}}))^2\right)^{1/2}$.

Figure 1 gives the mean and the quantiles of probability 5% and 95% of the RMSE computed from the 500 sets $(D_{2,i}, D_{1,i})_{i=1, \dots, 500}$ when the number of runs for the expensive code n_2 varies. In Figure 1, we can see that the errors converge to the same value when n_2 tends

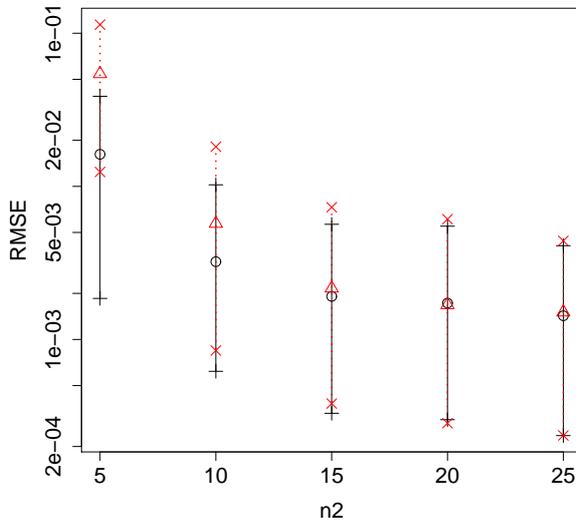


Figure 1: Comparison between kriging and co-kriging with $n_1 = 25$ runs for the cheap code (500 nested design sets have been randomly generated for each n_2). The circles represent the averaged RMSE of the co-kriging, the triangles represent the averaged RMSE of the kriging, the crosses represent the quantiles of probability 5% and 95% for the co-kriging RMSE and the times signs represent the quantiles of probability 5% and 95% of the kriging RMSE. Co-kriging predictions are better than the ordinary kriging ones for small n_2 and they converge to the same accuracy when n_2 tends to $n_1 = 25$.

to n_1 . Indeed, due to the Markov property given in Section 2.1, when $D_2 = D_1$, only the observations z_2 are taken into account. Furthermore, we can see that for small values of n_2 , it is worth considering the co-kriging model since its accuracy is significantly better than the one of the kriging model.

5.3 Nested space filling design

As presented in Section 2 we consider nested experimental design sets: $\forall t = 2, \dots, s \quad D_t \subseteq D_{t-1}$. Therefore, we have to adopt particular design strategies to uniformly spread the inputs for all D_t . A strategy based on Orthogonal array-based Latin hypercube for nested space-filling designs is proposed by [Qian et al., 2009].

We consider here another strategy for space-filling design, described in the following algorithm, which is very simple and not time-consuming. The number of points n_t for each design D_t is prescribed by the user, as well as the experimental design method applied to determine the coarsest grid D_s used for the most expensive code z_s (see [Fang et al., 2006] for a review of different methods).

ALGORITHM

build $D_s = \{x_j^{(s)}\}_{j=1, \dots, n_s}$ with the experimental design method prescribed by the user.

for $t = s$ to 2 **do**:

build design \tilde{D}_{t-1} with the experimental design method prescribed by the user.

for $i = 1$ to n_t **do**:

 find $\tilde{x}_j^{(t-1)} \in \tilde{D}_{t-1}$ the closest point from $x_i^{(t)} \in D_t$ where $j \in [1, n_{t-1}]$.

 remove $\tilde{x}_j^{(t-1)}$ from \tilde{D}_{t-1} .

end for

$D_{t-1} = \tilde{D}_{t-1} \cup D_t$.

end for

This strategy allows us to use any space-filling design method and it conserves the initial structure of the experimental design D_s of the most accurate code, contrarily to a strategy based on selection of subsets of an experimental design for the less accurate code as presented by [Kennedy and O'Hagan, 2000] and [Forrester et al., 2007]. We hence can ensure that D_s has excellent space-filling properties. Moreover, the experimental design D_{t-1} being equal to $\tilde{D}_{t-1} \cup D_t$, this method ensure the nested property.

In the presented application, we consider $n_2 = 5$ points for the expensive code $z_2(x)$ and $n_1 = 25$ points for the cheap one $z_1(x)$. We apply the previous algorithm to build D_2 and D_1 such that $D_2 \subset D_1$. For the experimental design set D_2 , we use a Latin-Hypercube-Sampling [Stein, 1987] optimized with respect to the S-optimality criterion which maximizes the mean distance from each design point to all the other points [Stocki, 2005]. Furthermore, the set D_1 is built using a maximum entropy design [Shewry and Wynn, 1987] optimized with the Fedorov-Mitchell exchange algorithm [Currin et al., 1991]. These algorithms are implemented in the library R lhs. The obtained nested designs are shown in Figure 2.

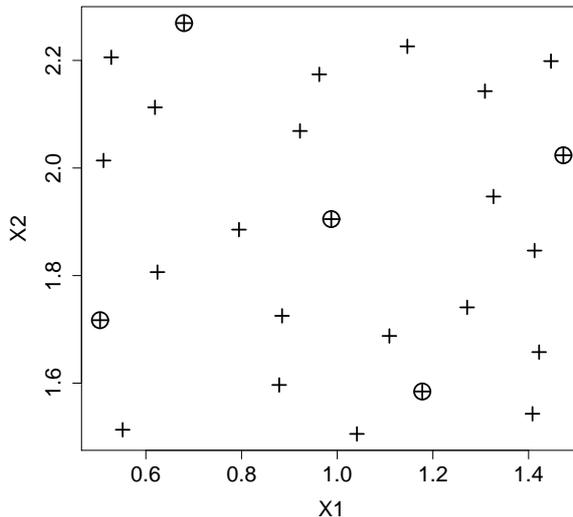


Figure 2: Nested experimental design sets for the hydrodynamic application. The crosses represent the $n_1 = 25$ points of the experimental design set D_1 of the cheap code and the circles represent the $n_2 = 5$ points of the experimental design set D_2 of the expensive code.

5.4 Multi-fidelity surrogate model for the dissipation factor *eps*

We build here a co-kriging model for the dissipation factor *eps*. The obtained model is compared to a kriging one. This first example is used to illustrate the efficiency of the co-kriging method compared to the kriging. It will also allow us to highlight the difference between the simple and the universal co-kriging.

We use the experimental design sets presented in Section 5.3. To validate and compare our models, the 175 simulations of the complex code uniformly spread on $[0.5, 1.5] \times [1.5, 2.3]$ are used. To build the different correlation matrices, we consider a tensorised matern- $\frac{5}{2}$ kernel (see [Rasmussen and Williams, 2006]):

$$r(x, x'; \theta_t) = r_{1d}(x_1, x'_1; \theta_{t,1})r_{1d}(x_2, x'_2; \theta_{t,2}) \quad (25)$$

with $x = (x_1, x_2) \in [0.5, 1.5] \times [1.5, 2.3]$, $\theta_{t,1}, \theta_{t,2} \in \mathbb{R}$ and:

$$r_{1d}(x_i, x'_i; \theta_{t,i}) = \left(1 + \sqrt{5} \frac{|x_i - x'_i|}{\theta_{t,i}} + \frac{5}{3} \frac{(x_i - x'_i)^2}{\theta_{t,i}^2} \right) \exp \left(-\sqrt{5} \frac{|x_i - x'_i|}{\theta_{t,i}} \right) \quad (26)$$

Then, we consider $g_1(x) = 1$, $f_2(x) = 1$, $f_1(x) = 1$ (see Section 2.1 and 2.2) and, using the concentrated maximum likelihood (see Section 5.1), we have the following estimations for the correlation hyper-parameters: $\hat{\theta}_1 = (0.69, 1.20)$ and $\hat{\theta}_2 = (0.27, 1.37)$.

According to the values of the hyper-parameter estimates, the co-kriging model are very smooth since the correlation lengths are large compared to the size of the input parameter space. Furthermore, the estimated correlation between the two codes is 82.64%, which shows that the amount of information contained in the cheap code is substantial.

Table 1 presents the results of the parameter estimations (see Section 2.3).

Trend coefficient	$\Sigma_t \nu_t$	Σ_t / σ_t^2
β_1	8.84	0.48
$\begin{pmatrix} \beta_{\rho_1} \\ \beta_2 \end{pmatrix}$	$\begin{pmatrix} 0.92 \\ 0.74 \end{pmatrix}$	$\begin{pmatrix} 1.98 & -18.13 \\ -18.13 & 165.82 \end{pmatrix}$
Variance coefficient	Q_t	$2\alpha_t$
σ_1^2	6.98	24
σ_2^2	0.06	3

Table 1: Application: hydrodynamic simulator. Parameter estimation results for the response *eps* (see equations (14) and (16)).

We see in Table 1 that the correlation between β_{ρ_1} and β_2 is important which highlights the importance of taking into account the correlation between these two coefficients for the parameter estimation. We also see that the adjustment parameter β_{ρ_1} is close to 1, which means that the two codes are highly correlated.

Figure 3 illustrates the contour plot of the kriging and co-kriging mean, we can see significant differences between the two surrogate models.

Table 2 compares the prediction accuracy of the co-kriging and the kriging models. The different coefficients are estimated with the 175 responses of the complex code on the test set:

- MaxAE: Maximal absolute value of the observed error.

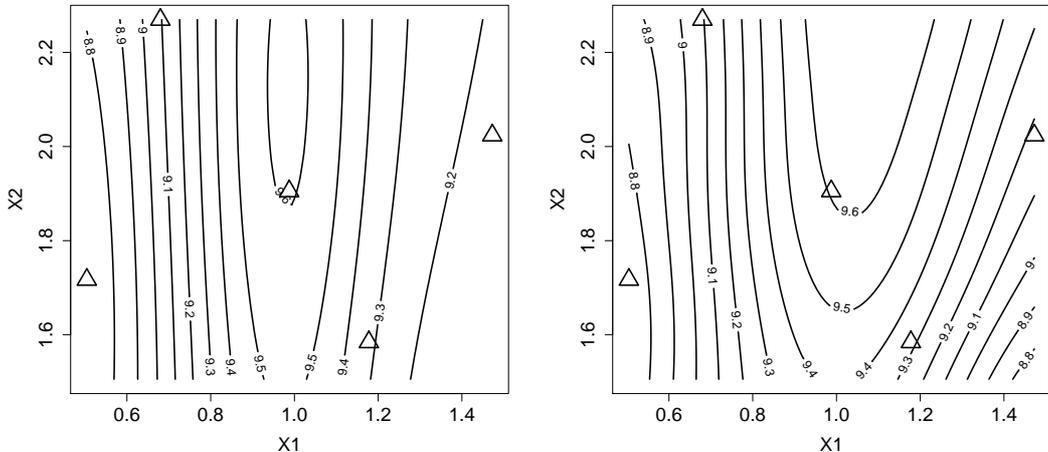


Figure 3: Contour plot of the kriging mean (on the left hand side) and the co-kriging mean (on the right hand side). The triangles represent the $n_2 = 25$ points of the experimental design set of the expensive code.

- RMSE : Root mean squared value of the observed error.
- $Q_2 = 1 - \|\mu_{Z_2}(D_{\text{test}}) - z_2(D_{\text{test}})\|^2 / \|\mu_{Z_2}(D_{\text{test}}) - \bar{z}_2\|^2$, with $\bar{z}_2 = (\sum_{i=1}^{n_2} z_2(x_i^{\text{test}})) / n_2$.
- RIMSE : Root of the average value of the kriging or co-kriging variance.

	Q_2	RMSE	MaxAE	RIMSE.
kriging	75.83%	0.133	0.49	0.110
co-kriging	98.01%	0.038	0.14	0.046

Table 2: Application: hydrodynamic simulator. Comparison between kriging and co-kriging. The co-kriging model provides predictions significantly better than the ones of the kriging model.

We can see that the difference of accuracy between the two models is important. Indeed, the one of the co-kriging model is significantly better. Furthermore, comparing the RMSE and the RIMSE estimations in Table 2, we see that we have a good estimation of the predictive distribution variances for the two models. We note that the predictive variance for the co-kriging is obtained with a simple co-kriging model. Therefore, it will be slightly larger in the universal co-kriging case. Indeed, by computing the universal co-kriging equations, we find $\text{RIMSE} = 0.058$.

We can compare the RMSE obtained with the test set with the RMSE obtained with a Leave-One-Out cross validation procedure (see Section 4). For this procedure, we test our model on $n_2 = 5$ validation sets obtained by removing one observation at a time. As presented in Section 4, we can either choose to remove the observations from z_2 or from z_2 and z_1 . The root mean squared error of the Leave-One-Out cross validation procedure obtained

by removing observations from z_2 is $\text{RMSE}_{z_2,LOO} = 4.80 \cdot 10^{-3}$ whereas the one obtained by removing observations from z_2 and z_1 is $\text{RMSE}_{z_1,z_2,LOO} = 0.10$. Comparing $\text{RMSE}_{z_2,LOO}$ and $\text{RMSE}_{z_1,z_2,LOO}$ to the RMSE obtained with the external test set, we see that the procedure which consists in removing points from z_2 and z_1 provides a better proxy for the generalization error. Indeed, $\text{RMSE}_{z_2,LOO}$ is a relevant proxy for the generalization error only at points where z_1 is available. Therefore, it underestimates the error at locations where $z_1(x)$ is unknown.

Figure 4 represents the mean and confidence intervals at plus or minus twice the standard deviation of the simple and universal co-krigings for points along the vertical line $x_1 = 0.99$ and the horizontal line $x_2 = 1.91$ ($x = (0.99, 1.91)$ corresponds to the coordinates of the point of D_2 in the center of the domain $[0.5, 1.5] \times [1.5, 2.3]$).

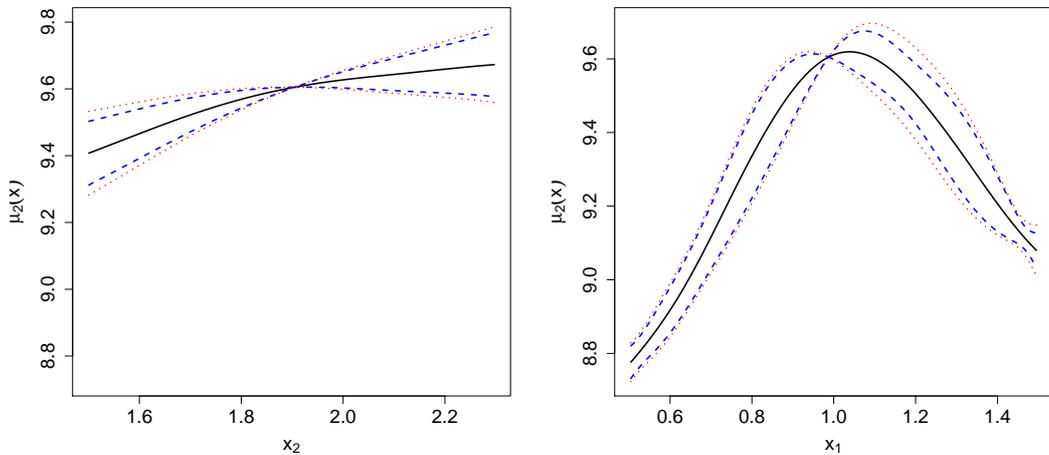


Figure 4: Mean and confidence intervals for the simple and the universal co-kriging. The figure on the left hand side represents the predictions along the vertical line $x_1 = 0.99$ and the figure on the right hand side represents the predictions along the horizontal line $x_2 = 1.91$. The solid black lines represent the mean of the two co-kriging models, the dashed lines represent the confidence interval at plus or minus twice the standard deviation of the simple co-kriging and the dotted lines represent the same confidence intervals for the universal co-kriging.

In Figure 4 on the right hand side, we see a necked point around the coordinates $x_1 = 1.5$ since, in the direction of x_2 , the hyper-parameters of correlation for $Z_1(x)$ and $\delta_2(x)$ are large ($\theta_{1,2} = 1.20$ and $\theta_{2,2} = 1.37$) and a point of D_2 have almost the same coordinate.

5.5 Multi-fidelity surrogate model for the mixture characteristic length L_c

In this section, we build a co-kriging model for the mixture characteristic length L_c . The aim of this example is to highlight that it could be worth having an adjustment coefficient ρ_1 depending on x . We use the same training and test sets as in the previous section and we consider a tensorised matern- $\frac{5}{2}$ kernel (25). Let us consider the two following cases:

- Case 1: $g_1(x) = 1$, $f_2(x) = 1$ and $f_1(x) = 1$
- Case 2: $g_1^T(x) = (1 \ x_1)$, $f_2(x) = 1$ and $f_1(x) = 1$

We have the following hyper-parameter maximum likelihood estimates for the two cases

- Case 1: $\hat{\theta}_1 = (0.52, 1.09)$ and $\hat{\theta}_2 = (0.03, 0.02)$
- Case 2: $\hat{\theta}_1 = (0.52, 1.09)$ and $\hat{\theta}_2 = (0.14, 1.37)$

The estimation of $\hat{\theta}_1$ is identical in the two cases since it does not depend on ρ_1 and it is estimated with the same observations. Furthermore, we see an important difference between the estimates of $\hat{\theta}_2$. Indeed, they are larger in the Case 2 than in the Case 1 which suppose that the model is smoother in the Case 2. Table 3 presents the estimations of β_1 and σ_1^2 for the two cases (see Section 2.3).

Trend coefficient	$\Sigma_1 \nu_1$	Σ_1 / σ_1^2
β_1	1.26	0.97
Variance coefficient	Q_1	$2\alpha_1$
σ_1^2	15.62	24

Table 3: Application: hydrodynamic simulator. Estimations of β_1 and σ_1^2 for the response L_c (see equations (14) and (16)).

Then, Table 4 presents the estimations of β_2 , β_{ρ_1} and σ_2^2 for the Case 1, i.e. when ρ_1 is constant (see Section 2.3).

Trend coefficient	$\Sigma_2 \nu_2$	Σ_2 / σ_2^2
$\begin{pmatrix} \beta_{\rho_1} \\ \beta_2 \end{pmatrix}$	$\begin{pmatrix} 1.49 \\ -0.26 \end{pmatrix}$	$\begin{pmatrix} 0.83 & -0.79 \\ -0.79 & 0.95 \end{pmatrix}$
Variance coefficient	Q_2	$2\alpha_2$
σ_2^2	0.01	3

Table 4: Application: hydrodynamic simulator. Estimations of β_2 , β_{ρ_1} and σ_2^2 for the Case 1, i.e. when ρ_1 is constant, for the response L_c (see equations (14) and (16)).

Finally, Table 5 presents the estimations of β_2 , β_{ρ_1} and σ_2^2 for the Case 2, i.e. when ρ_1 depends on x (see Section 2.3).

Trend coefficient	$\Sigma_2 \nu_2$	Σ_2 / σ_2^2
$\begin{pmatrix} \beta_{\rho_1} \\ \beta_2 \end{pmatrix}$	$\begin{pmatrix} 1.66 \\ -0.48 \\ -0.04 \end{pmatrix}$	$\begin{pmatrix} 2.34 & -3.50 & 0.44 \\ -3.50 & 9.18 & -3.67 \\ 0.44 & -3.67 & 2.60 \end{pmatrix}$
Variance coefficient	Q_2	$2\alpha_2$
σ_2^2	$3.24 \cdot 10^{-4}$	2

Table 5: Application: hydrodynamic simulator. Estimations of β_2 , β_{ρ_1} and σ_2^2 for the Case 2, i.e. when ρ_1 depends on x , for the response L_c (see equations (14) and (16)).

We see in Table 4 that the adjustment coefficient is around 1.5 which indicates that the magnitude of the expensive code is slightly more important than the one of the cheap code. Furthermore, we see in Table 5 that if we consider an adjustment coefficient which linearly depends on x_1 (i.e. with $g_1^T(x) = (1 \ x_1)$), the constant part of β_{ρ_1} is more important (it is around 1.66) and there is a negative slope in the direction x_1 (it is around -0.48). Since $x \in [0.5, 1.5]$, the averaged value of ρ_1 is 1.18 and goes from 1.42 at $x_1 = 0.5$ to 0.94 at $x_1 = 1.5$. We see also a significant difference between the two case for the variance estimation. Indeed, the variance estimate in the Case 1 (see Table 4) is much more important than the one in the Case 2 (see Table 5). This could mean that we learn better in the Case 2 than in the Case 1.

Figure 5 illustrates the contour plot of the two co-kriging models, i.e. when ρ_1 is constant and when ρ depends on x .

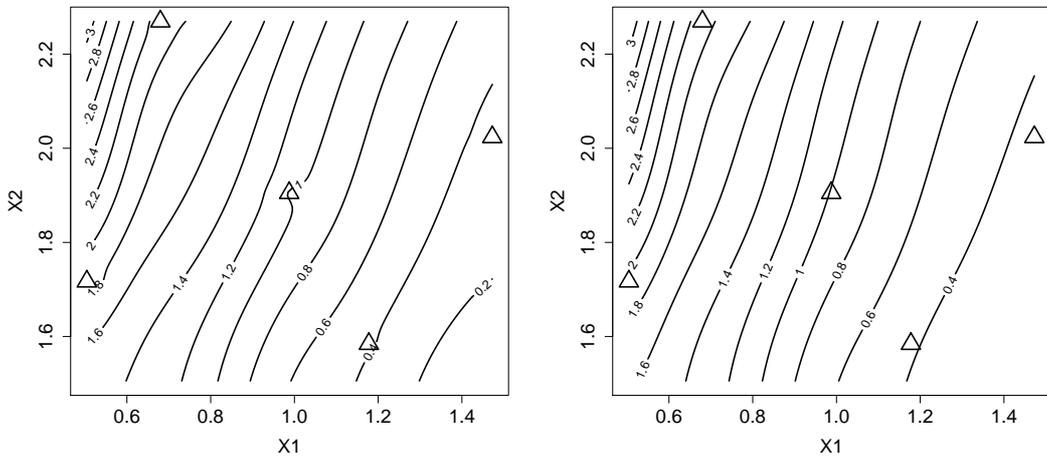


Figure 5: Contour plot of the co-kriging mean when ρ_1 is constant (on the left hand side) and when ρ_1 is depends on x (of the right hand side). The triangles represent the $n_2 = 5$ points of the experimental design set of the expensive code.

Furthermore, Table 6 compares the prediction accuracy of the co-kriging in the two cases. The precision is computed on the test set of 175 observations.

	RMSE	MaxAE
Case 1	$7.26 \cdot 10^{-3}$	0.23
case 2	$1.53 \cdot 10^{-3}$	0.16

Table 6: Application: hydrodynamic simulator. Comparison between co-kriging when ρ_1 is constant (Case 1) and co-kriging when ρ_1 depends on x (Case 2). The Case 2 provides predictions better than the Case 1, it is hence worthwhile to consider an adjustment coefficient not constant.

We see that the co-kriging model in Case 2 is clearly better than the one in Case 1. Therefore, we illustrate in this application that it can be worth considering an adjustment

coefficient not constant contrary to the model presented in [Kennedy and O’Hagan, 2000] and [Forrester et al., 2007].

6 Conclusion

We have presented in this paper a recursive formulation for a multi-fidelity co-kriging model. This model allows us to build surrogate models using data from experiments of different levels of fidelity.

The strength of the suggested approach is that it considerably reduces the complexity of the co-kriging model while it preserves its predictive efficiency. Therefore, the proposed method is competitive regarding the Bayes linear approach in which the principal strength is a low computational cost but with a low predictive efficiency. Furthermore, one of the most important consequences of the recursive formulation is that the construction of the surrogate model is equivalent to build s independent krigings. Consequently, we can naturally adapt results of kriging to the proposed co-kriging model.

In this paper, we first prove that our model is equivalent to another very popular model in terms of predictive distributions whereas it reduces its complexity. Then, we present a Bayesian estimation of the model parameters which provides closed form expressions for the parameters of the posterior distributions. We note that, from these posterior distributions, we can deduce the maximum likelihood estimates of the parameters. Then, thanks to the joint distributions of the parameters and the recursive formulation, we can deduce closed form formulas for the mean and covariance of the posterior predictive distribution. Due to their similarities with the universal kriging equations, we call these formulas the universal co-kriging equations. Finally, we present closed form expressions for the cross-validation equations of the co-kriging surrogate model. These expressions reduce considerably the complexity of the cross-validation procedure and are derived from the one of kriging model that we have extended.

The suggested model has been successfully applied to a hydrodynamic code. We also present in this application a practical way to design the experiments of the multi-fidelity model and we show that it is worth using this co-kriging model instead of a kriging model.

7 Acknowledgements

The author particularly thanks Professor Josselin Garnier for supervising his work and for his fruitful guidance. He is also grateful to Dr. Claire Cannaméla for providing the data for the application and for her interesting discussions and helpful comments.

A Proofs

A.1 Proof of Proposition 1

Let us consider the co-kriging mean of the model (1) presented in [Kennedy and O’Hagan, 2000] for a t -level co-kriging with $t = 2, \dots, s$:

$$m_{Z_t}(x) = h'_t(x)^T \beta^{(t)} + t_t(x)^T V_t^{-1} (z^{(t)} - H_t \beta^{(t)})$$

where $\beta^{(t)} = (\beta_1^T, \dots, \beta_t^T)^T$, $z^{(t)} = (z_1^T, \dots, z_t^T)^T$ and $h'_t(x)^T$ is defined in equation (6). We have:

$$\begin{aligned} h'_t(x)^T \beta^{(t)} &= \rho_{t-1}(x) \left(\left(\prod_{i=1}^{t-2} \rho_i(x) \right) f_1^T(x), \left(\prod_{i=2}^{t-2} \rho_i(x) \right) f_2^T(x), \dots, f_{t-1}^T(x) \right) \beta^{(t-1)} + f_t^T(x) \beta_t \\ &= \rho_{t-1}(x) h'_{t-1}(x)^T \beta^{(t-1)} + f_t^T(x) \beta_t \end{aligned}$$

Then, from equations (7) and (8), we have the following equality:

$$\begin{aligned} t_t(x)^T V_t^{-1} z^{(t)} &= \rho_{t-1}(x) t_{t-1}(x)^T V_{t-1}^{-1} z^{(t-1)} - (\rho_{t-1}^T(D_t)) \odot (r_t^T(x) R_t^{-1} z_{t-1}(D_t)) \\ &\quad + r_t^T(x) R_t^{-1} z_t \end{aligned}$$

and with equation (6):

$$t_t(x)^T V_t^{-1} H_t \beta^{(t)} = \rho_{t-1}(x) t_{t-1}(x)^T V_{t-1}^{-1} H_{t-1} \beta^{(t-1)} + r_t^T(x) R_t^{-1} F_t(D_t) \beta_t$$

where \odot stands for the element by element matrix product. We hence obtain the recursive relation:

$$m_{Z_t}(x) = \rho_{t-1}(x) m_{Z_{t-1}}(x) + f_t^T(x) \beta_t + r_t^T(x) R_t^{-1} [z_t - \rho_{t-1}(D_t) \odot z_{t-1}(D_t) - F_t(D_t) \beta_t]$$

The co-kriging mean of the model (9) satisfies the same recursive relation (6), and we have $m_{Z_1}(x) = \mu_{Z_1}(x)$. This proves the first equality of Proposition 1:

$$\mu_{Z_s}(x) = m_{Z_s}(x)$$

We follow the same guideline for the co-kriging covariance:

$$s_{Z_t}^2(x, x') = v_{Z_t}^2(x, x') - t_t^T(x) V_t^{-1} t_t(x')$$

where $v_{Z_t}^2(x, x')$ is the covariance between $Z_t(x)$ and $Z_t(x')$ and $s_{Z_t}^2(x, x')$ is the covariance function of the conditioned Gaussian process $[Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \beta, \beta_\rho, \sigma^2]$ for the model (1). From equation (8), we can deduce the following equality:

$$\sigma_{Z_t}^2(x, x') = \rho_{t-1}(x) \rho_{t-1}(x') v_{Z_{t-1}}^2(x, x') + v_t^2(x, x')$$

where $\sigma_{Z_t}^2(x, x')$ is the covariance function of the conditioned Gaussian process $[Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \beta_t, \beta_{\rho_{t-1}}, \sigma_t^2]$ of the recursive model (9). Then, from equation (7) and (8), we have:

$$t_t^T(x) V_t^{-1} t_t(x') = \rho_{t-1}(x) \rho_{t-1}(x') t_{t-1}^T(x) V_{t-1}^{-1} t_{t-1}(x') + \sigma_t^2 r_t^T(x) R_t^{-1} r_t(x')$$

Finally we can deduce the following equality:

$$s_{Z_t}^2(x, x') = \rho_{t-1}(x) \rho_{t-1}(x') \left(v_{Z_{t-1}}^2(x, x') - t_{t-1}^T(x) V_{t-1}^{-1} t_{t-1}(x') \right) + \sigma_t^2 (1 - r_t^T(x) R_t^{-1} r_t(x'))$$

which is equivalent to:

$$s_{Z_t}^2(x, x') = \rho_{t-1}(x) \rho_{t-1}(x') s_{Z_{t-1}}^2(x, x') + \sigma_t^2 (1 - r_t^T(x) R_t^{-1} r_t(x'))$$

This is the same recursive relation as the one satisfies by the co-kriging covariance $\sigma_{Z_t}^2(x, x')$ of the model (9) (see equation (12)). Since $s_{Z_1}^2(x, x') = \sigma_{Z_1}^2(x, x')$, we have :

$$\sigma_{Z_s}^2(x, x') = s_{Z_s}^2(x, x')$$

This equality with $x = x'$ proves the second equality of Proposition 1. \square

A.2 Proof of Proposition 2

Noting that the mean of the predictive distribution in equation (10) do not depend on σ_t^2 and thanks to the law of total expectation, we have the following equality:

$$\mathbb{E} \left[Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)} \right] = \mathbb{E} \left[\mathbb{E} \left[Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \sigma_t^2, \beta_t, \beta_{\rho_{t-1}} \right] \middle| \mathcal{Z}^{(t)} = z^{(t)} \right]$$

From the equations (11) and (14), we directly deduce the equation (17). Then, we have the following equality:

$$\text{var} \left(\mu_{Z_t}(x) \middle| z^{(t)}, \sigma_t^2 \right) = (h_t^T(x) - r_t(x)^T R_t^{-1} H_t) \Sigma_t (h_t^T(x) - r_t(x)^T R_t^{-1} H_t)^T \quad (27)$$

The law of total variance states that:

$$\begin{aligned} \text{var}(Z_t(x) | z^{(t)}, \sigma_t^2) &= \mathbb{E} \left[\text{var}(Z_t(x) | z^{(t)}, \beta_t, \beta_{\rho_{t-1}}, \sigma_t^2) \middle| z^{(t)}, \sigma_t^2 \right] \\ &\quad + \text{var} \left(\mathbb{E} \left[Z_t(x) | z^{(t)}, \beta_t, \beta_{\rho_{t-1}}, \sigma_t^2 \right] \middle| z^{(t)}, \sigma_t^2 \right) \end{aligned}$$

Thus, from equations (11), (17) and (27), we obtain:

$$\begin{aligned} \text{var}(Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \sigma_t^2) &= \hat{\rho}_t^2(x) \text{var}(Z_{t-1}(x) | \mathcal{Z}^{(t-1)} = z^{(t-1)}, \sigma_t^2) + \sigma_t^2 (1 - r_t^T(x) R_t^{-1} r_t^T(x)) \\ &\quad + (h_t^T - r_t^T(x) R_t^{-1} H_t) \Sigma_t (h_t^T - r_t^T(x) R_t^{-1} H_t)^T \end{aligned} \quad (28)$$

Again using the law of total variance and the independence between $\mathbb{E} [Z_t(x) | \mathcal{Z}^{(t)} = z^{(t)}, \beta_t, \beta_{\rho_{t-1}}]$ and σ_t^2 , we have:

$$\text{var}(Z_t(x) | z^{(t)}) = \mathbb{E}_{\sigma_t^2} \left[\text{var}(Z_t(x) | z^{(t)}, \sigma_t^2) \right] \quad (29)$$

We obtain the equation (18) from equation (16) by noting that the mean of an inverse Gamma distribution $\mathcal{IG}(a, b)$ is $b/(a-1)$ \square

A.3 Proof of Proposition 3

Let us consider that ξ_s is the index of the k last points of D_s . We denote by D_{test} these points. First we consider the variance and the trend parameters as fixed, i.e. $\sigma_{t, -\xi_t}^2 = \frac{Q_t}{2(a_t-1)}$ and $\lambda_{t, -\xi_t} = \Sigma_t \nu_t$, and $\mathcal{V}_s = 0$, i.e. we are in the simple co-kriging case. Thanks to the block-wise inversion formula, we have the following equality:

$$R_s^{-1} = \begin{pmatrix} A & B \\ B^T & Q^{-1} \end{pmatrix} \quad (30)$$

with $A = [R_s^{-1}]_{[-\xi_s, -\xi_s]} + [R_s^{-1}]_{[-\xi_s, -\xi_s]} [R_s^{-1}]_{[-\xi_s, \xi_s]} Q^{-1} [R_s^{-1}]_{[\xi_s, -\xi_s]} [R_s^{-1}]_{[-\xi_s, -\xi_s]}$,
 $B = - [R_s^{-1}]_{[-\xi_s, -\xi_s]} [R_s^{-1}]_{[-\xi_s, \xi_s]} Q^{-1}$ and:

$$Q = [R_s^{-1}]_{[\xi_s, \xi_s]} - [R_s^{-1}]_{[\xi_s, -\xi_s]} \left([R_s^{-1}]_{[-\xi_s, -\xi_s]} \right)^{-1} [R_s^{-1}]_{[-\xi_s, \xi_s]} \quad (31)$$

We note that $\frac{Q_s}{2(a_s-1)} Q = \frac{Q_t}{2(a_t-1)} \left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1}$ represents the covariance matrix of the points in D_{test} with respect to the covariance kernel of a Gaussian process of kernel $\frac{Q_s}{2(a_s-1)} r_s(x, x')$

(which is the one of $\delta_s(x)$) conditioned by the points $D_s \setminus D_{\text{test}}$. Therefore, from the previous remark and the equation (12), we can deduce the equation (21).

Furthermore, we have the following equality:

$$\begin{aligned} \left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1} [R_s^{-1} (z_s - H_s \lambda_{s, -\xi_s})]_{[\xi_s]} &= z_s(D_{\text{test}}) - h_s^T(D_{\text{test}}) \Sigma_s \nu_s \\ &- [R_s^{-1}]_{[-\xi_s, \xi_s]} \left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1} \\ &\times (z_{s-1}(D_s \setminus D_{\text{test}}) - [H_s^T]_{[-\xi_s]} \Sigma_s \nu_s) \end{aligned} \quad (32)$$

From this equation and equation (11), we can directly deduce the equation (19) with $\varepsilon_{Z_s, \xi_s} = z_s(D_{\text{test}}) - \mu_{Z_s}(D_{\text{test}})$.

Then, we suppose the trend and the variance parameters as unknown and we have to re-estimate them when we remove the observations. Thanks to the parameter estimations presented in Section 2.3, we can deduce that the estimates of $\sigma_{t, -\xi_t}^2$ and $\lambda_{t, -\xi_t}$ when we remove observations of index ξ_t are given by the following equations:

$$\lambda_{s, -\xi_s} ([H_s^T]_{-\xi_s} K_s [H_s]_{-\xi_s}) = [H_s^T]_{-\xi_s} K_s z_s(D_{\text{test}}) \quad (33)$$

and:

$$\sigma_{s, -\xi_s}^2 = \frac{(z_s(D_{\text{test}}) - [H_s]_{-\xi_s} \lambda_{s, -\xi_s})^T K_s (z_s(D_{\text{test}}) - [H_s]_{-\xi_s} \lambda_{s, -\xi_s})}{n_s - p_s - q_{s-1} - n_{\text{train}}} \quad (34)$$

with $K_s = \left([R_s]_{[-\xi_s, -\xi_s]} \right)^{-1}$.

From the equality (30), we can deduce that $K_s = A - BQB^T$ from which we obtain the equation (20). Finally, to obtain the cross-validation equations for the universal co-kriging, we just have to estimate the following quantity (see equation (18)):

$$\left(h_s^T(D_{\text{test}})^T - [R_s^{-1}]_{[-\xi_s, \xi_s]} K_s [H_s]_{-\xi_s} \right) \Sigma_s \left(h_s^T(D_{\text{test}})^T - [R_s^{-1}]_{[-\xi_s, \xi_s]} K_s [H_s]_{-\xi_s} \right)^T \quad (35)$$

with $\Sigma_s = \left([H_s^T]_{-\xi_s} K_s [H_s]_{-\xi_s} \right)^{-1}$. The following equality:

$$\left(h_s^T(D_{\text{test}})^T - [R_s^{-1}]_{[-\xi_s, \xi_s]} K_s [H_s]_{-\xi_s} \right) = \left(\left([R_s^{-1}]_{[\xi_s, \xi_s]} \right)^{-1} [R_s^{-1} H_s]_{[\xi_s]} \right) \quad (36)$$

allows us to obtain the equation (23) and completes the proof. \square

References

- [Craig et al., 1998] Craig, P. S., Goldstein, M., Seheult, A. H., and Smith, J. A. (1998). Constructing partial prior specifications for models of complex physical systems. *Applied Statistics*, 47:37–53.
- [Cumming and Goldstein, 2009] Cumming, J. A. and Goldstein, M. (2009). Small sample bayesian designs for complex high-dimensional models based on information gained using fast approximations. *Technometrics*, 51:377–388.
- [Currin et al., 1991] Currin, C., Mitchell, T., Morris, M., and Ylvisaker, D. (1991). Bayesian prediction of deterministic functions with applications to the design and analysis of computer experiments. *American Statistical Association*, 86:953–963.

- [Dubrule, 1983] Dubrule, O. (1983). Cross validation of kriging in a unique neighborhood. *Mathematical Geology*, 15:687–699.
- [Fang et al., 2006] Fang, K.-T., Li, R., and Sudjianto, A. (2006). *Design and Modeling for Computer Experiments*. Computer Science and Data Analysis Series, London.
- [Forrester et al., 2007] Forrester, A. I. J., Sobester, A., and Keane, A. J. (2007). Multi-fidelity optimization via surrogate modelling. *Proc. R. Soc. A*, 463:3251–3269.
- [Goldstein and Wooff, 2007] Goldstein, M. and Wooff, D. A. (2007). *Bayes Linear Statistics: Theory and Methods*. Chichester, England: Wiley.
- [Grégoire et al., 2005] Grégoire, O., Souffland, D., and Serge, G. (2005). A second order turbulence model for gaseous mixtures induced by richtmyer-meshkov instability. *Journal of Turbulence*, 6:1–20.
- [Jeffreys, 1961] Jeffreys, H. (1961). *Theory of Probability*. Oxford University Press, London.
- [Kennedy and O’Hagan, 2000] Kennedy, M. C. and O’Hagan, A. (2000). Predicting the output from a complex computer code when fast approximations are available. *Biometrika*, 87:1–13.
- [Krige, 1951] Krige, D. G. (1951). A statistical approach to some basic mine valuation problems on the witwatersrand. *Technometrics*, 52:119–139.
- [Matheron, 1963] Matheron, G. (1963). Principles of geostatistics. *Economic Geology*, 58:1246–1266.
- [Matheron, 1969] Matheron, G. (1969). *Le krigeage Universel*. Ecole des Mines de Paris, Paris.
- [Patterson and Thompson, 1971] Patterson, H. and Thompson (1971). Recovery of interblock information when block sizes are unequal. *Biometrika*, 58:545–554.
- [Qian et al., 2009] Qian, P. Z. G., Ai, M., and Wu, C. F. J. (2009). Construction of nested space-filling designs. *The Annals of Statistics*, 37:3616–3643.
- [Qian and Wu, 2008] Qian, P. Z. G. and Wu, C. F. J. (2008). Bayesian hierarchical modeling for integrating low-accuracy and high-accuracy experiments. *Technometrics*, 50:192–204.
- [Rasmussen and Williams, 2006] Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press, Cambridge.
- [Sacks et al., 1989] Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989). Design and analysis of computer experiments. *Statistical Science*, 4:409–423.
- [Santner et al., 2003] Santner, T. J., Williams, B. J., and Notz, W. I. (2003). *The Design and Analysis of Computer Experiments*. Springer, New York.
- [Shewry and Wynn, 1987] Shewry, M. C. and Wynn, H. P. (1987). Maximum entropy sampling. *Journal of Applied Statistics*, 14:165–170.

- [Stein, 1987] Stein, M. L. (1987). Large sample properties of simulations using latin hypercube sampling. *Technometrics*, 29:143–151.
- [Stein, 1999] Stein, M. L. (1999). *Interpolation of Spatial Data*. Springer Series in Statistics, New York.
- [Stocki, 2005] Stocki, R. (2005). A method to improve design reliability using optimal latin hypercube sampling. *Computer Assisted Mechanics and Engineering Sciences*, 12:87–105.
- [Zhang and Wang, 2009] Zhang, H. and Wang, Y. (2009). Kriging and cross-validation for massive spatial data. *Environmetrics*, 21:290–304.