



HAL
open science

Gaussian process regression of two nested computer codes

Sophie Marque-Pucheu

► **To cite this version:**

Sophie Marque-Pucheu. Gaussian process regression of two nested computer codes. Mathematics [math]. Université Paris-Diderot - Paris VII, 2018. English. NNT: . tel-02092072v2

HAL Id: tel-02092072

<https://hal.science/tel-02092072v2>

Submitted on 7 May 2019 (v2), last revised 9 Dec 2019 (v4)

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.

Thèse de doctorat
de l'Université Sorbonne Paris Cité
Préparée à l'Université Paris Diderot
Ecole doctorale n°386 Mathématiques Paris Centre
Laboratoire de Probabilités, Statistique et Modélisation

Gaussian process regression of two nested computer codes

Par Sophie Marque-Pucheu

Thèse de doctorat de Mathématiques Appliquées

Présentée et soutenue publiquement à Paris le 10 octobre 2018 devant le jury suivant

Examinatrice	Fischer Aurélie	Maître de conférences	Université Paris Diderot
Directeur de thèse	Garnier Josselin	Professeur	École Polytechnique
Examinatrice	Marrel Amandine	Ingénieur de recherche	CEA
Rapporteur	Monod Hervé	Directeur de recherche	INRA
Président du jury	Nouy Anthony	Professeur	École Centrale Nantes
Examineur	Perrin Guillaume	Ingénieur de recherche	CEA

D'après les rapports de

Monod Hervé	Directeur de recherche	INRA
Marzouk Youssef	Professor	MIT

Introduction

Remerciements

Mes premiers remerciements s'adressent naturellement aux deux encadrants qui m'ont accompagnée lors de ces trois années de thèse. Josselin Garnier, mon directeur de thèse, a suivi mes travaux de manière attentive. Sa très grande culture scientifique et sa réactivité ont été d'une aide précieuse. Guillaume Perrin, mon encadrant au CEA, a également fait preuve d'une très grande implication dans le suivi de cette thèse. Je le remercie pour sa pédagogie et sa confiance.

Je remercie également les deux rapporteurs de cette thèse : Hervé Monod et Youssef Marzouk. Thank you, Mr. Marzouk, for your careful reading of the manuscript. Merci également à M. Monod pour sa lecture attentive du manuscrit, ainsi que pour avoir été membre du jury.

Je tiens aussi à remercier Aurélie Fischer, Amandine Marrel et Anthony Nouy d'avoir accepté de faire partie de mon jury de thèse.

Enfin, je salue tous les membres de l'équipe incertitudes du CEA, ainsi que les doctorants du bâtiment B.

Je remercie également tous les personnels administratifs du CEA et l'Université Paris Diderot qui ont facilité d'une manière ou d'une autre mon quotidien.

Mes derniers remerciements vont à ma famille, en particulier mes parents qui m'ont transmis le goût d'apprendre. Last but not least, je remercie enfin mon conjoint pour son formidable soutien.

Résumé

Cette thèse traite de la métamodélisation (ou émulation) par processus gaussien de deux codes couplés. Le terme « deux codes couplés » désigne ici un système de deux codes chaînés : la sortie du premier code est une des entrées du second code.

Les deux codes sont coûteux. Afin de réaliser une analyse de sensibilité de la sortie du code couplé, on cherche à construire un métamodèle de cette sortie à partir d'un faible nombre d'observations. Trois types d'observations du système existent : celles de la chaîne complète, celles du premier code uniquement, celles du second code uniquement. Le métamodèle obtenu doit être précis dans les zones les plus probables de l'espace d'entrée.

Les métamodèles sont obtenus par krigeage universel, avec une approche bayésienne.

Dans un premier temps, le cas sans information intermédiaire, avec sortie scalaire, est traité. Une méthode innovante de définition de la fonction de la moyenne du processus gaussien, basée sur le couplage de deux polynômes, est proposée. Ensuite le cas avec information intermédiaire est traité. Un prédicteur basé sur le couplage des prédicteurs gaussiens associés aux deux codes est proposé. Des méthodes pour évaluer rapidement la moyenne et la variance du prédicteur obtenu sont proposées. Les résultats obtenus pour le cas scalaire sont ensuite étendus au cas où les deux codes sont à sortie de grande dimension. Pour ce faire, une méthode de réduction de dimension efficace de la variable intermédiaire de grande dimension est proposée pour faciliter la régression par processus gaussien du deuxième code. Les méthodes proposées sont appliquées sur des exemples numériques.

Mots-clés

Codes numériques emboîtés, codes couplés, codes chaînés, régression par processus gaussien, métamodélisation, variable fonctionnelle, réduction de dimension, Stepwise Uncertainty Reduction, plans d'expériences séquentiels.

Abstract

This thesis deals with the Gaussian process regression of two nested codes. The term "nested codes" refers to a system of two chained computer codes: the output of the first code is one of the inputs of the second code.

The two codes are computationally expensive. In order to perform a sensitivity analysis, we aim at emulating the output of the nested code from a small number of observations.

Three types of observations of the system exist: those of the chained code, those of the first code only and those of the second code only. The surrogate model has to be accurate on the most likely regions of the input domain of the nested code.

In this work, the surrogate models are constructed using the Universal Kriging framework, with a Bayesian approach.

First, the case when there is no information about the intermediary variable (the output of the first code) is addressed. An innovative parametrization of the mean function of the Gaussian process modeling the nested code is proposed. It is based on the coupling of two polynomials. Then, the case with intermediary observations is addressed. A stochastic predictor based on the coupling of the predictors associated with the two codes is proposed. Methods aiming at computing quickly the mean and the variance of this predictor are proposed. Finally, the methods obtained for the case of codes with scalar outputs are extended to the case of codes with high dimensional vectorial outputs.

We propose an efficient dimension reduction method of the high dimensional vectorial input of the second code in order to facilitate the Gaussian process regression of this code.

All the proposed methods are applied to numerical examples.

Keywords

Nested computer codes, Gaussian process regression, surrogate modeling, functional variable, dimension reduction, Stepwise Uncertainty Reduction, sequential designs.

Résumé long en français

Cette thèse présente de nouveaux développements pour la métamodélisation de codes coûteux chaînés, où la sortie du premier code est une des entrées du code suivant. Cette configuration et sa généralisation à plus que deux codes sont fréquemment rencontrées en pratique. Mais la construction de métamodèles adaptés à cette configuration a été peu étudiée jusqu'ici.

Ce manuscrit contient trois contributions nouvelles par rapport à l'état de l'art, détaillées dans les chapitres 3 à 5. La première contribution concerne la régression par processus gaussien avec une fonction de moyenne définie par un polynôme. Une nouvelle méthode de définition de la tendance polynomiale, basée sur la composition de deux polynômes, est proposée. Dans ce cas de figure, la variable intermédiaire entre les deux codes n'est pas connue.

La seconde contribution suppose la connaissance de la variable intermédiaire et traite de l'enrichissement du plan d'expériences en vue de la régression par processus gaussien de la sortie de la chaîne de deux codes. Le choix d'une nouvelle observation soulève plusieurs questions. Tout d'abord pour un code donné, il faut choisir les variables d'entrée de la nouvelle observation. Ensuite, comme il y a deux codes, la question se pose également (si cela est possible) de choisir auquel des deux codes ajouter une nouvelle observation.

La troisième contribution traite le cas de deux codes à sortie de très grande dimension (par exemple des fonctions du temps). Dans cette configuration, le second code a une sortie, mais également une entrée fonctionnelle. Une méthode de réduction de dimension de l'entrée fonctionnelle adaptée à ce cas est alors proposée. Les critères d'enrichissement proposés précédemment sont combinés avec cette méthode de réduction de dimension afin de les étendre au cas de deux codes à sortie fonctionnelle. Les méthodes proposées sont ensuite appliquées à un cas test industriel modélisant l'explosion d'une charge dans une cuve sphérique. Ce cas test est associé à un couplage entre un code de détonique et un code de dynamique des structures. Les paragraphes qui suivent présentent plus en détails la structure du manuscrit.

Le premier chapitre passe en revue l'état de l'art concernant la métamodélisation d'un unique code à entrée et sortie de faibles dimensions. Une brève présentation de la régression linéaire et du chaos polynomial est faite, ainsi que de méthodes de régularisation comme LASSO ou LARS. Le reste du chapitre est dédié à la régression par processus gaussien (GP) ou krigeage. Après un rappel des bases de la régression par processus gaussien, comme le choix de la fonction de covariance, le krigeage universel dans un cadre bayésien est présenté. Ensuite, les critères pour plans d'expériences pour la régression par processus gaussien et l'optimisation bayésienne sont passés en revue. Le chapitre se conclut sur une brève partie concernant l'analyse de sensibilité, en particulier les méthodes basées sur une décomposition de la variance (indices de Sobol).

Le deuxième chapitre passe en revue les méthodes pour la régression par processus gaussien d'un code à entrée et/ou sortie définie comme une fonction discrétisée du temps. L'attention se concentre ici sur la réduction de la dimension de l'entrée ou de la sortie. Concernant la réduction de la dimension de l'entrée, certaines méthodes ne prennent en compte que l'entrée fonctionnelle, tandis que d'autres ont pour objectif la réduction de la dimension de l'entrée de manière adaptée à la sortie. Ces dernières sont tout particulièrement adaptées pour le système chaîné considéré dans ce travail. Concernant la sortie fonctionnelle, deux approches sont possibles. La première consiste à projeter la sortie fonctionnelle sur une base de dimension réduite. La seconde repose sur l'utilisation d'une covariance tensorisée, où l'indice de la sortie

fonctionnelle (comme par exemple le temps) est considéré comme une des entrées du modèle.

Le troisième chapitre contient la première contribution de cette thèse : la construction d'une fonction de moyenne du processus gaussien par couplage de deux polynômes. Cette approche intègre l'information que l'on a a priori sur la structure chaînée des deux codes, mais sans observations ni connaissance de la structure de la variable intermédiaire. Dans ce cas, la configuration est proche d'une régression par processus gaussien classique, avec des observations des entrées et sortie de la chaîne de codes. La spécificité de la méthode repose sur l'utilisation de l'information que l'on a sur cette structure chaînée. La définition de la fonction de moyenne comprend une première étape de composition de deux polynômes, puis une seconde étape de linéarisation de cette composition. Cette linéarisation permet de limiter l'impact d'une erreur d'estimation des paramètres de chacun des deux polynômes. Ensuite le prédicteur de la sortie de la chaîne de code est construit en utilisant le krigeage universel dans un cadre bayésien. Par ailleurs, la structure proposée pour la tendance polynomiale offre une grande flexibilité, puisque les ordres totaux de chacun des deux polynômes, mais aussi la dimension de la sortie du premier polynôme, peuvent être optimisés. Cependant, cette flexibilité nécessite la résolution d'un problème d'optimisation complexe car non convexe. Une approche heuristique, basée sur une minimisation alternée par rapport aux variables, est proposée pour résoudre ce problème d'optimisation. Par ailleurs, un critère basé sur l'erreur Leave One Out (LOO) est utilisé pour caractériser la performance de prédiction du prédicteur gaussien. Ce critère est utilisé pour choisir la combinaison de valeurs la plus performante pour les ordres totaux des deux polynômes et la dimension de la sortie du premier polynôme.

Le quatrième chapitre contient la deuxième contribution de cette thèse : la métamodélisation de deux codes chaînés lorsque des observations de la variable intermédiaire sont disponibles. Le prédicteur proposé est basé sur un couplage de prédicteurs gaussiens de chacun des deux codes. Le chapitre propose en particulier deux critères d'enrichissement du plan d'expériences. Ces critères reposent sur une minimisation de la variance de prédiction intégrée (IMSE). La variance de prédiction doit donc être évaluée en un très grand nombre de points. Le premier critère correspond au cas où les deux codes ne peuvent pas être appelés de manière séparée. Le second correspond au cas où les codes peuvent être lancés de manière séparée. Dans ce cas, on peut choisir lequel des deux codes appeler, en retenant celui qui maximise la réduction de la variance de prédiction intégrée par unité de temps de calcul pour une évaluation du code. Une difficulté majeure liée à cette approche tient au fait que le couplage de deux prédicteurs gaussiens n'est pas gaussien. La variance de prédiction doit donc être évaluée en utilisant des méthodes de quadrature ou Monte Carlo. Afin de résoudre ces difficultés numériques, deux méthodes pour une évaluation rapide de la variance de prédiction sont proposées. Dans le premier cas, si le processus gaussien associé au second code a une fonction de covariance gaussienne et une tendance polynomiale, alors la variance peut être évaluée de manière analytique. Dans le cas où ces conditions ne sont pas valables, une autre approche reposant sur la linéarisation du couplage des deux prédicteurs peut être utilisée. Les méthodes proposées sont ensuite appliquées sur deux exemples numériques : un premier analytique et un second portant sur la trajectoire balistique d'un projectile conique. Les résultats obtenus montrent l'intérêt de prendre en compte les observations de la variable intermédiaire et de pouvoir appeler de manière séparée chacun des deux codes.

Le cinquième chapitre contient les contributions finales de cette thèse et concerne la métamodélisation par processus gaussien de deux codes chaînés à sortie fonctionnelle (de très grande dimension). La contribution majeure de ce chapitre est une méthode de réduction de l'entrée fonctionnelle d'un modèle linéaire, qui est adaptée à la sortie de ce modèle linéaire.

Cette méthode de réduction de dimension est combinée à une approximation de la sortie du second code, qui est linéaire par rapport à l'entrée fonctionnelle du second code (qui est également la sortie du premier code). Le modèle linéaire proposé est en fait un filtre causal, paramétré par un petit nombre de variables qui peuvent être estimées à partir d'un faible nombre d'observations.

Cette combinaison d'une approximation linéaire et d'une réduction de dimension adaptée à ce modèle linéaire permet de réduire la dimension de l'entrée fonctionnelle du second code de manière adaptée à la prédiction de la sortie de ce code.

Grâce à cette réduction de dimension, chacun des deux codes peut être associé à un processus gaussien avec un vecteur d'entrées de faible dimension. Deux prédicteurs gaussiens sont obtenus en utilisant une covariance tensorisée pour prendre en compte le caractère multidimensionnel des sorties des fonctions considérées. Les prédicteurs sont ensuite couplés et le couplage est linéarisé. Ceci permet d'obtenir un prédicteur gaussien de la sortie fonctionnelle de la chaîne de deux codes. La moyenne et la variance du prédicteur peuvent alors être évaluées de manière analytique, et donc très rapide. Les critères d'enrichissement proposés dans le chapitre précédent sont ensuite adaptés au cas de deux codes couplés à sortie fonctionnelle. Enfin, les méthodes proposées sont mises en application sur le cas test industriel qui a motivé cette thèse, à savoir le couplage d'un code de détonique avec un code de dynamique des structures. Les sorties de chacun des codes sont des fonctions discrétisées du temps. Les résultats obtenus montrent l'intérêt de prendre en compte les observations de la variable intermédiaire, par rapport à une simple régression par processus gaussien de la sortie de la chaîne de codes en fonction des entrées.

Contents

Introduction	i
Notations	xiii
I State of the art for the surrogate modeling of computer codes	1
1 Surrogate modeling of a single code with scalar inputs and output	5
1.1 Linear regression	5
1.2 Polynomial Chaos Expansion	6
1.3 Methods for the selection of the regressors of a linear model	8
1.3.1 Stepwise and all-subsets regressions	8
1.3.2 Ridge regression	8
1.3.3 LASSO	9
1.3.4 Forward stagewise regression	9
1.3.5 Least Angle Regression	9
1.3.6 Dantzig selector	11
1.3.7 Conclusions	11
1.4 Gaussian process regression or Kriging	11
1.4.1 Gaussian processes	11
1.4.2 Ordinary, simple and universal Kriging	18
1.4.3 Estimation of a parametric covariance function	20
1.5 Design of experiments	22
1.5.1 Space-filling designs	22
1.5.2 Criterion-based designs	24
1.5.3 Gaussian processes for pointwise global optimization	26
1.6 Sensitivity analysis	26
2 Gaussian process regression of a code with a functional input or output	29
2.1 Dimension reduction of a functional variable	29
2.1.1 Dimension reduction adapted to the functional variable only	30
2.1.2 Dimension reduction adapted to the functional variable and a dependent variable	31
2.2 Gaussian process prediction of a computer code with a functional output	33
2.2.1 Projection of the functional output on a basis	34
2.2.2 Gaussian process regression of the whole functional output	35

II	Contributions	37
3	Nested polynomial trends for the improvement of Gaussian predictors	39
3.1	Introduction	39
3.2	Gaussian process predictors	41
3.2.1	General framework	41
3.2.2	Choice of the covariance function	42
3.2.3	Choice of the mean function	42
3.3	Nested polynomial trends for Gaussian process predictors	43
3.3.1	Nested polynomial representations	43
3.3.2	Coupling nested representations and Gaussian processes	46
3.3.3	Linearization of the nested polynomial trend	47
3.3.4	Error evaluation	48
3.3.5	Convergence analysis	50
3.4	Applications	50
3.4.1	$d = 1$	51
3.4.2	$d > 1$	54
3.4.3	Relevance of the LOO error	55
3.5	Conclusions	57
4	Gaussian process regression of two nested codes with scalar output	59
4.1	Introduction	59
4.2	Surrogate modeling for two nested computer codes	60
4.2.1	General framework	60
4.2.2	Gaussian process-based surrogate models	61
4.2.3	Sequential designs for the improvement of Gaussian process predictors	64
4.3	Fast computation of the variance of the predictor of the nested code	65
4.3.1	Explicit derivation of the two first statistical moments of the predictor	66
4.3.2	Linearized approach	67
4.4	Applications	68
4.4.1	Characteristics of the examples	69
4.4.2	Prediction performance for a given set of observations	72
4.4.3	Performances of the sequential designs	74
4.5	Conclusions	78
4.6	Proofs	79
4.6.1	Proof of Proposition 4.2.1	79
4.6.2	Proof of Lemma 4.3.1	79
4.6.3	Proof of Lemma 4.3.2	80
4.6.4	Proof of Proposition 4.3.1	81
4.6.5	Proof of Proposition 4.3.2	85
4.6.6	Proof of Corollary 4.3.3	85
	Conclusions	89

Context

Surrogate modeling for the sensitivity analysis of two nested computer codes

This thesis is motivated by an application case. This application case is the coupling of two computationally costly computer codes. The first code is a detonation code and the second code is a structural dynamics code. The two codes have functional (i.e. high dimensional vectorial) outputs and the functional output of the first code is one of the inputs of the second code.

If we aim at performing design and certification studies of such a system, the evaluation of the output of the system at a large number of input points is often necessary. This is especially true when methods like sensitivity analysis, risk analysis or optimization are performed.

In this work we aim at performing a sensitivity analysis of the system mentioned above. Given the computational cost of the two codes, the first objective is to build an emulator, or a surrogate model, of the output of the two nested codes. This surrogate model will be constructed from a small set of observations of the two codes. The number of observations cannot be very high because of the computational costs of the codes.

As the role of simulation is increasing, the surrogate modeling of high-cost codes generates growing interest. However, the existing methods are generally applied to a single code or consider a system of codes as a single code.

In this work, the framework of the Gaussian process regression for the surrogate modeling of computer codes is considered. In this framework, the output of a code is considered to be the realization of a Gaussian process. The framework used for the Gaussian process regression is the Universal Kriging framework and a Bayesian approach is utilized. If some not very restrictive assumptions on the prior distribution of the Gaussian process are fulfilled, a Gaussian predictor of the code can be obtained by computing the posterior distribution of the Gaussian process given the observations of the code output.

Moreover, the existing methods for the surrogate modeling of codes generally consider the case of codes with low dimensional vectorial inputs. If a code has a functional input, the dimension of the functional input is often reduced thanks to a projection. The choice of the optimal method of dimension reduction of the functional input for the surrogate modeling of the output remains a research topic.

Contributions of the thesis

This thesis makes contributions to the surrogate modeling of two nested codes with scalar or functional outputs. These contributions aim at solving the following difficulties of the studied system:

- there are two codes,
- the codes are coupled by a functional intermediary variable,
- the second code has a functional input.

First, the case of two nested codes with scalar outputs is investigated. The considered system is then:

$$\begin{array}{ccc} & \mathbf{x}_2 & \\ & \searrow & \\ \mathbf{x}_1 & \rightarrow & y_1(\mathbf{x}_1) \nearrow \\ & & y_{\text{nest}}(\mathbf{x}_{\text{nest}}) := y_2(y_1(\mathbf{x}_1), \mathbf{x}_2), \end{array} \quad (0.0.1)$$

with $\mathbf{x}_1 \in \mathbb{R}^{d_1}$ and $\mathbf{x}_2 \in \mathbb{R}^{d_2}$ the low dimensional vectorial inputs of the two codes, $y_1 \in \mathbb{R}$ and $y_2 \in \mathbb{R}$ the output of the two codes, and d_1 and d_2 two integers.

In a first step, the case where there are no observations of the intermediary variable $y_1(\mathbf{x}_1)$ is considered. An innovative parametrization of the mean function of the Gaussian process is proposed. This parametrization is based on the coupling of polynomials and enables to improve the prediction accuracy compared to a classical constant or polynomial mean function.

Then the case where observations of the intermediary variable are available is considered. A stochastic predictor of the nested code is obtained by coupling the Gaussian predictors of the two codes. Such an approach enables to take into account all the types of observations: observations of the nested code, of the first code only and of the second code only. The predictor is non-Gaussian but its moments can be computed using Monte Carlo methods. Then we define sequential design criteria which aim at improving the prediction accuracy of the proposed predictor. The criteria are based on a reduction of the integrated prediction variance because the predictor has to be accurate on the most probable areas of the input domain for the sensitivity analysis. Finally, two adaptations of the proposed predictor are developed in order to evaluate the prediction variance and thus the proposed sequential design criteria quickly. The first adaptation is called "analytic" and the second one "linearized". They both enable to compute the mean and the variance of the proposed predictor in closed forms. The "linearized" method leads also to a Gaussian predictor of the nested code. Moreover, the interest of taking into account the intermediary observations is shown.

Finally, the case of two nested code with functional outputs is investigated. The considered system is then:

$$\begin{array}{ccc} & \mathbf{x}_2 & \\ & \searrow & \\ \mathbf{x}_1 & \rightarrow & \mathbf{y}_1(\mathbf{x}_1) \end{array} \quad \begin{array}{c} \searrow \\ \nearrow \end{array} \quad \mathbf{y}_{\text{nest}}(\mathbf{x}_{\text{nest}}) := \mathbf{y}_2(\mathbf{y}_1(\mathbf{x}_1), \mathbf{x}_2), \quad (0.0.2)$$

with $\mathbf{y}_1 \in \mathbb{R}^{N_t}$ and $\mathbf{y}_2 \in \mathbb{R}^{N_t}$ the output of the two codes when they are functional, $N_t \gg 1$ denoting the number of discretization steps of the functional outputs.

The second code has a functional input and the existing methods of Gaussian process regression generally consider low dimensional vectorial inputs. The Gaussian process regression of the second code requires therefore the reduction of the dimension of this functional input. We propose a dimension reduction of the functional input of a code which is suited for the prediction of the functional output of this code. This dimension reduction method is based on a two-step approach. First, the output of the second code is approximated by a linear causal filter. This linear model has a sparse structure, which is defined by only N_t variables. These variables can be estimated from a small set of observations of the functional input and output of the second code. The second step is the use of a proposed projection basis which is adapted to a linear model. The combination of these two steps enables to obtain a dimension reduction of the functional input of the second code, which:

- is adapted to the output of this code
- can be estimated from a small set of observations,
- does not require the knowledge of the derivatives of the output of the code,

Once the dimension of the functional intermediary variable has been efficiently reduced, the previously defined linearized method is adapted to the case of two nested codes with functional outputs. A Gaussian predictor of the functional output of the nested code, with analytic mean and variance, is obtained. Finally, the previously defined sequential design criteria are adapted to the case of two nested codes with functional outputs.

Outline of the manuscript

The thesis has two parts.

Part I provides a review of the state of the art for the surrogate modeling of computer codes.

In Chapter 1, we review methods for the surrogate modeling of a single code with low dimensional vectorial inputs and a scalar output.

Section 1.1 describes the surrogate modeling of a single code by Linear Regression.

Section 1.2 focuses on the surrogate modeling of a code by Polynomial Chaos Expansion.

Section 1.3 reviews the existing methods for the selection of the regressors in the framework of Linear Regression.

Section 1.4 provides a review of the Gaussian process regression framework for the surrogate modeling of a single code with low dimensional vectorial inputs and a scalar output.

Section 1.5 presents a review of the design of experiments for an accurate surrogate model on the whole input domain of a code with low dimensional vectorial inputs and a scalar output.

Section 1.6 focuses on the sensitivity analysis of the output of a code, or a quantity associated with it, with respect to the inputs of the code.

In Chapter 2, we review methods for the surrogate modeling of a single code with a functional output, low dimensional vectorial inputs and possibly a functional input.

Section 2.1 is devoted to the existing methods for the dimension reduction of a functional variable.

Section 2.2 reviews the existing methods for the Gaussian process regression of a code with low dimensional vectorial inputs and a functional output.

Part II details our contributions to the construction of a surrogate model of two nested codes with scalar or functional outputs.

In Chapter 3, we focus on the case where the two codes have scalar outputs and no observations of the intermediary variable are available. We propose to define the mean function of the Gaussian process modeling the nested code as a coupling of two polynomials. This parametrization is based on the coupling of two polynomials. We show how this parametrization can improve the prediction accuracy of the Gaussian predictor compared to the case where the mean function is defined by polynomials.

In Chapter 4 we focus on the case where the two codes have scalar outputs and observations of the intermediary variable are available. We propose a stochastic predictor of the nested code based on the coupling of the Gaussian predictors of the two codes. This stochastic predictor is non-Gaussian but its mean and variance can be evaluated using Monte Carlo methods. This predictor can take into account all the possible observations: those of the nested code, those of the first code and those of the second code. Then sequential design criteria are proposed. These design criteria aim at improving the prediction accuracy on the whole input domain of the nested code. One of the criteria can also take into account the difference of computational costs between the two codes. Finally, we propose two adaptations of the previously proposed predictor of the nested code in order to accelerate the computation of the mean and the variance of the predictor. They both enable to compute the prediction mean and variance in closed forms. In addition, the proposed linearized predictor of the nested code enables to obtain a Gaussian predictor of the nested code with conditioned mean and variance functions in closed forms.

The application of the proposed methods to numerical examples shows the interest of taking into account the intermediary observations.

In Chapter ?? we focus on the case of the coupling of two codes with functional outputs. We first propose an efficient dimension reduction of the functional input of the second code. This dimension reduction is based on a linear projection of the functional input of the second code. The proposed projection basis can be estimated from a small set of observations of the second code and does not require the knowledge of the derivatives of the code.

We also extend the linearized predictor of the nested code proposed in Chapter 4 to the case of two nested codes with scalar output. This extension relies on the dimension reduction of the functional output and a tensorized structure of the Gaussian process modeling the code. By tensorized structure we mean a separation between the index of the output and the inputs. The sequential design criteria are also adapted to the case of two nested codes with functional outputs.

The proposed methods are applied to numerical examples. The results show again the interest of taking appropriately into account the intermediary observations.

The predictor obtained at the end of the sequential enrichment of the initial design is used in order to perform a sensitivity analysis of a scalar quantity of interest based on the functional output of the nested code.

Notations

Ordinal variables

n	number of observations
d	dimension of an input variable
p	number of functions of a basis of function in the case of Universal Kriging
N_t	dimension of the time-varying output of a code
$\text{card}(\mathbb{A})$	number of elements of the set \mathbb{A}

Matrix, vectors and scalar

x	a scalar
\mathbf{x}	a vector
x_i or $(\mathbf{x})_i$	the i -th entry of the vector \mathbf{x}
\mathbf{X}	a matrix
$(\mathbf{X})_{ij}$	the entry at line i and row j of the matrix \mathbf{X}
$(\mathbf{X})_{.i}$	the vector of the entries of the i -th column of the matrix \mathbf{X}
$(\mathbf{X})_i$	the vector of the entries of the i -th row of the matrix \mathbf{X}
\mathbf{X}^T	transpose of the matrix \mathbf{X}
$\text{diag}(\mathbf{x})$	diagonal matrix with diagonal \mathbf{x}
$\text{diag}(\mathbf{X})$	vector corresponding to the diagonal of the matrix \mathbf{X}
$\text{Tr}(\mathbf{X})$	trace of the matrix \mathbf{X}
$\text{cov}(\mathbf{x}, \mathbf{y})$	covariance between \mathbf{x} and \mathbf{y}

Probabilistic notations

$\stackrel{d}{=}$	equality in distribution
$\mathbb{E}[\cdot]$	Mean of a random quantity
$\mathbb{V}[\cdot]$	Variance of a random quantity
$\mathcal{N}(\mathbf{m}, \mathbf{K})$	multivariate normal distribution with mean \mathbf{m} and covariance matrix \mathbf{K}
$\text{GP}(m(\cdot), C(\cdot, \cdot))$	one-dimensional Gaussian process with mean function m and covariance function C
$\text{GP}(\mathbf{m}(\cdot), \mathbf{C}(\cdot, \cdot))$	multidimensional Gaussian process with vector-valued mean function \mathbf{m} and matrix-valued covariance function \mathbf{C}

Norms and scalar products

$(\cdot, \cdot)_{\mathbb{X}}$	scalar product in the space of square integrable real-valued functions on \mathbb{X} , such that $(y, z)_{\mathbb{X}} := \int_{\mathbb{X}} y(\mathbf{x})z(\mathbf{x})d\mathbf{x}$
$\ \cdot\ _{\mathbb{X}}$	norm in the space of square integrable real-valued functions on \mathbb{X} , such that $\ y\ _{\mathbb{X}}^2 := (y, y)_{\mathbb{X}}$
$\ \cdot\ _F$	Frobenius norm
$\ \cdot\ _1$	L_1 norm, such that $\ \mathbf{x}\ _1 = \sum_{i=1}^d x_i $
$\ \cdot\ $	L_2 norm, such that $\ \mathbf{x}\ _2 = \sqrt{\sum_{i=1}^d x_i^2}$

Part I

State of the art for the surrogate modeling of computer codes

The role of simulation for the design and the certification of complex systems is increasing. However, methods like uncertainty propagation, sensitivity analysis or optimization require the evaluation of the output of the code at a huge number of input points. If the computational cost of the computer code is high, and only a small number of observations of its output is available, the use of a surrogate model is necessary. In this part we review some existing methods for the surrogate modeling of computer codes.

This part includes two chapters. The first one is devoted to the surrogate modeling of a computer code with scalar (i.e. low dimensional vectorial) inputs and output. The second one focuses on the surrogate modeling with Gaussian process regression of a code with functional (i.e. high dimensional vectorial) input and/or output.

Chapter 1

Surrogate modeling of a single code with scalar inputs and output

In this chapter we consider a model of the form $\mathbf{x} \mapsto y(\mathbf{x})$, $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^d$, d a positive integer, and $\mu_{\mathbb{X}}$ is a probability measure on the space comprising \mathbb{X} and a σ -algebra over \mathbb{X} . The following sections detail the state of the art for the surrogate modeling of y from a set of n observations of the input and the output of the code. These observations are denoted by:

$$\mathbf{X}^{\text{obs}} = \begin{pmatrix} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{x}^{(n)} \end{pmatrix}, \quad (1.0.1)$$

and

$$\mathbf{y}^{\text{obs}} = \left(y^{(1)} = y(\mathbf{x}^{(1)}), \dots, y^{(n)} = y(\mathbf{x}^{(n)}) \right), \quad (1.0.2)$$

where \mathbf{X}^{obs} is a $(n \times d)$ -dimensional matrix and \mathbf{y}^{obs} is a n -dimensional vector.

The first section is devoted to linear regression. The second one deals with the use of Polynomial Chaos Expansion as a surrogate model. The third one focuses on the methods for the selection of regressors in regression models. The fourth one presents the Gaussian process regression for the surrogate modeling of a computer code. Finally, the last section reviews some existing designs of experiments which are adapted for the acquisition of knowledge of the computer code or the sequential improvement of a surrogate model.

1.1 Linear regression

Generalized additive models are a very common tool for the emulation of a response surface [Hastie and Tibshirani, 1990]. It is the projection of the output y on a basis of functions h_i , $1 \leq i \leq p$, p a positive integer, of the inputs \mathbf{x} . The emulator can be written in the form:

$$\hat{y}(\mathbf{x}) = \mathbf{h}(\mathbf{x})^T \boldsymbol{\beta}, \quad (1.1.1)$$

where $\mathbf{h}(\mathbf{x})$ and $\boldsymbol{\beta}$ are in \mathbb{R}^p . The functions of the basis can be polynomials, with Polynomial Chaos Expansion as a particular case, wavelets, trigonometric functions...

Note that simple linear regression can be regarded as a particular case of the generalized additive models, with a basis of functions comprising only the covariates: $\mathbf{h}(\mathbf{x}) = \mathbf{x}$.

The regression coefficients $\boldsymbol{\beta}$ can be estimated from a set of n observations of the inputs and the output of the code \mathbf{X}^{obs} and \mathbf{y}^{obs} through the minimization of the quadratic loss function:

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \mathbf{h}(\mathbf{x}^{(i)})^T \boldsymbol{\beta} \right)^2. \quad (1.1.2)$$

If we denote:

$$\mathbf{H} = \begin{pmatrix} \mathbf{h}(\mathbf{x}^{(1)})^T \\ \vdots \\ \mathbf{h}(\mathbf{x}^{(n)})^T \end{pmatrix}, \quad (1.1.3)$$

then the least squares estimate of the regression coefficients can be written:

$$\hat{\boldsymbol{\beta}} = \mathbf{H}^+ \mathbf{y}^{\text{obs}}, \quad (1.1.4)$$

where \mathbf{H}^+ is the pseudo-inverse of \mathbf{H} . If $n \geq p$ and \mathbf{H} is of rank p , then $\mathbf{H}^T \mathbf{H}$ is invertible and $\mathbf{H}^+ = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$. By definition, \mathbf{H} is a $(n \times p)$ -dimensional matrix.

However, matrix $(\mathbf{H}^T \mathbf{H})$ is not always invertible. The number of observations can be smaller than the number of regression coefficients ($p \leq n$) or the functions of the basis can be correlated according to the probability measure $\mu_{\mathbb{X}}$, which means that the columns of \mathbf{H} are correlated, thus reducing the rank of matrix \mathbf{H} .

The matrix $(\mathbf{H}^T \mathbf{H})$ is more likely to be inverted if the basis functions are decorrelated with respect to the probability measure $\mu_{\mathbb{X}}$ of the inputs, as performed with Polynomial Chaos Expansion. Another possible approach is the use of a regularization term for the inversion of the matrix, or the selection of the most influencing regressors. The two following sections detail these two approaches.

1.2 Polynomial Chaos Expansion

Polynomial Chaos expansion can be used to emulate a model response y with inputs \mathbf{x} . Besides, the probability measure $\mu_{\mathbb{X}}$ associated with \mathbf{x} is a product measure. Therefore, the components of the input vector are independent. It has been applied by Ghanem and Spanos [1990] to stochastic finite elements methods. Polynomial Chaos expansion can be seen as the projection of the model output y on a polynomial basis which depends on the distribution of the model inputs \mathbf{x} . The polynomials are orthonormal with respect to the distribution of \mathbf{x} . The model response can therefore be expanded as:

$$y(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in \mathbb{N}^d} \beta_{\boldsymbol{\alpha}} \Phi_{\boldsymbol{\alpha}}(\mathbf{x}), \quad (1.2.1)$$

with $\beta_{\boldsymbol{\alpha}} \in \mathbb{R}$ and $\Phi_{\boldsymbol{\alpha}}$ orthonormal multidimensional polynomials, which means:

$$\int_{\mathbb{X}} \Phi_{\boldsymbol{\alpha}}(\mathbf{x}) \Phi_{\boldsymbol{\gamma}}(\mathbf{x}) d\mu_{\mathbb{X}}(\mathbf{x}) = \delta_{\boldsymbol{\alpha}\boldsymbol{\gamma}}, \quad (1.2.2)$$

with $\delta_{\boldsymbol{\alpha}\boldsymbol{\gamma}}$ denoting the Kronecker delta.

In practice, the expansion of Eq. (1.2.1) can be truncated in order to obtain a surrogate model of the model response. If we denote by $\mathbb{A} \subset \mathbb{N}^d$ the truncated set of indices, by $\boldsymbol{\beta}_{\mathbb{A}}$ the vector gathering the $\beta_{\boldsymbol{\alpha}}$, $\boldsymbol{\alpha} \in \mathbb{A}$ and by $\Phi_{\mathbb{A}}$ the vector gathering the selected polynomials, this surrogate model is defined as:

$$\hat{y}(\mathbf{x}) = \Phi_{\mathbb{A}}(\mathbf{x})^T \boldsymbol{\beta}_{\mathbb{A}}. \quad (1.2.3)$$

Note that the truncation is generally defined by an upper bound r on the total order of the polynomials, which means $\mathbb{A} = \{\boldsymbol{\alpha} \in \mathbb{N}^d, \|\boldsymbol{\alpha}\|_1 \leq r\}$. The total order r can be chosen adaptively according to a target precision, with an estimation of the error thanks to a cross-validation criterion [Blatman and Sudret, 2010, 2011].

A coefficient $\beta_{\boldsymbol{\alpha}}$ is defined as the projection of the model response on function $\Phi_{\boldsymbol{\alpha}}$:

$$\beta_{\boldsymbol{\alpha}} = \int_{\mathbb{X}} y(\mathbf{x}) \Phi_{\boldsymbol{\alpha}}(\mathbf{x}) d\mu_{\mathbb{X}}(\mathbf{x}). \quad (1.2.4)$$

Distribution	Density	Orthonormal basis
Uniform	$\frac{1}{2} \mathbb{1}_{[-1,1]}(x)$	$\frac{P_k(x)}{\sqrt{2k+1}}$, with P_k Legendre polynomial
Gaussian	$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$	$\frac{H_k(x)}{\sqrt{k!}}$, with H_k Hermite polynomial
Gamma	$\frac{x^a}{\Gamma(a+1)} \exp(-x) \mathbb{1}_{x>0}$	$\frac{L_k(x)}{\Gamma(k+a+1)}$, with L_k Laguerre polynomial

Table 1.1: Classical univariate polynomial families used for Polynomial Chaos Expansion.

The integral can be estimated using Monte-Carlo methods, quadrature rules [Ghiocel and Ghanem, 2002] or stochastic collocation methods [Xiu, 2009].

The coefficients can also be estimated by least squares regression [Blatman and Sudret, 2010, 2011] from a set of n observations:

$$\hat{\beta}_{\mathbb{A}} = \operatorname{argmin}_{\beta_{\mathbb{A}} \in \mathbb{R}^{\operatorname{card}(\mathbb{A})}} \sum_{i=1}^n \left(y^{(i)} - \Phi_{\mathbb{A}}(\mathbf{x}^{(i)})^T \beta_{\mathbb{A}} \right)^2. \quad (1.2.5)$$

Note that if the observations are drawn according to the distribution of the inputs, the meta-model will be more accurate in the high-probability regions of the input domain.

The usual one-dimensional polynomial families used for Polynomial Chaos Expansion, which are chosen according to the distribution of the one-dimensional variable x , are given in Table 1.1.

Furthermore, the inputs can be transformed using an isoprobabilistic transformation, such as the Nataf or the Rosenblatt transformations [Nataf, 1962; Rosenblatt, 1952; Lebrun and Dutfoy, 2009]. Such transformations map \mathbf{x} to a d -dimensional standard Gaussian variable $\boldsymbol{\xi}$ (i.e. d independent standard Gaussian variables). Then a Polynomial Chaos Expansion can be performed using Hermite polynomials [Blatman and Sudret, 2011]. The expansion becomes:

$$y(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in \mathbb{N}^d} \beta_{\boldsymbol{\alpha}} \mathcal{H}_{\boldsymbol{\alpha}}(T(\mathbf{x})), \quad (1.2.6)$$

where $\mathcal{H}_{\boldsymbol{\alpha}} = \prod_{i=1}^d H_{\alpha_i}$ and

$$\beta_{\boldsymbol{\alpha}} = \int_{T(\mathbb{X})} y(T^{-1}(\boldsymbol{\xi})) \mathcal{H}_{\boldsymbol{\alpha}}(\boldsymbol{\xi}) \prod_{i=1}^d \varphi(\xi_i) d\boldsymbol{\xi}. \quad (1.2.7)$$

Here, $T : \mathbf{x} \mapsto \boldsymbol{\xi}$ is the isoprobabilistic transformation and T^{-1} its inverse, $\mathcal{H}_{\boldsymbol{\alpha}}$ are Hermite polynomials, and φ the standard univariate Gaussian probability density function.

Thanks to this isoprobabilistic transformation, the Polynomial Chaos Expansion of a computer code with dependent inputs can be performed.

1.3 Methods for the selection of the regressors of a linear model

In this section we review the existing methods for the selection of the most influential regressors for linear regression or Polynomial Chaos Expansion. The methods are presented in the chronological order of their appearance. Two approaches can be distinguished: the first one selects the regressors which are the most influential. The second one minimizes the coefficients associated with the least influential regressors.

1.3.1 Stepwise and all-subsets regressions

Stepwise regression aims at selecting the regressors which improve the prediction accuracy the most. There are three main approaches to perform this selection: forward selection, backward elimination and bidirectional elimination.

In the forward method, the set of the selected regressors is empty at the initial step. Then, at each step, one adds the regressor which best improves the prediction accuracy of the regression model. The addition continues until a stopping criterion is reached.

On the contrary, with the backward elimination, a huge number of regressors are selected at the initial step. Then the regressors which contribute the least to the prediction accuracy are removed step by step from the regression model.

Efroymson [1960] introduced an approach combining forward selection and backward elimination. At each step of the forward selection, the interest of removing one of the previous selected regressors is studied.

However, stepwise regression is known as being greedy and quite unstable [Hesterberg et al., 2008].

In parallel, all-subsets regression has been introduced by Furnival and Wilson [1974]. It relies on the evaluation of the accuracy of all the regression models based on all the subsets of the set of regressors. Even though exhaustive, this approach can be computationally expensive, especially when the number of regressors is high.

1.3.2 Ridge regression

Introduced by Hoerl and Kennard [1970], ridge regression is based on a penalization of the coefficients of the regressors. This penalization can be seen as a regularization of the regression problem. The coefficients obtained with the ridge regression are the solutions of the following optimization problem:

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \mathbf{h}(\mathbf{x}^{(i)})^T \boldsymbol{\beta} \right)^2 + \delta \|\boldsymbol{\beta}\|_2^2, \quad (1.3.1)$$

with δ a non-negative real-valued constant.

This leads to the normal equation:

$$(\mathbf{H}^T \mathbf{H} + \delta \mathbf{I}_p) \hat{\boldsymbol{\beta}}^{\text{ridge}} = \mathbf{H}^T \mathbf{y}^{\text{obs}}. \quad (1.3.2)$$

Practically, the optimal value of δ can be estimated thanks to a Cross validation criterion. The absolute value of the coefficients decreases as δ increases. When $\delta = 0$, the result is the same as the one of ordinary least squares. If $\delta > 0$ then the matrix $(\mathbf{H}^T \mathbf{H} + \delta \mathbf{I}_p)$ is positive definite and thus invertible.

The ridge regression can be seen as a particular case of the Tikhonov regularization [Tikhonov

and Arsenin, 1977], which is defined as follows:

$$\widehat{\boldsymbol{\beta}}^{\text{Tikhonov}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \mathbf{h}(\mathbf{x}^{(i)})^T \boldsymbol{\beta} \right)^2 + \|\boldsymbol{\Gamma} \boldsymbol{\beta}\|^2, \quad (1.3.3)$$

with $\boldsymbol{\Gamma}$ a $d \times d$ -dimensional matrix.

If $\boldsymbol{\Gamma}^T \boldsymbol{\Gamma}$ is positive definite, this problem has the following explicit solution:

$$\widehat{\boldsymbol{\beta}}^{\text{Tikhonov}} = (\mathbf{H}^T \mathbf{H} + \boldsymbol{\Gamma}^T \boldsymbol{\Gamma})^{-1} \mathbf{H}^T \mathbf{y}^{\text{obs}}. \quad (1.3.4)$$

Note that if $\boldsymbol{\Gamma}$ is defined such that $\boldsymbol{\Gamma}^T \boldsymbol{\Gamma}$ is positive definite, then the matrix $\mathbf{H}^T \mathbf{H} + \boldsymbol{\Gamma}^T \boldsymbol{\Gamma}$ is an invertible matrix.

1.3.3 LASSO

The Least Absolute Shrinkage and Selection Operator (LASSO) method has been introduced by Tibshirani [1989]. It relies on a L_1 -penalization of the estimation of $\boldsymbol{\beta}$, which can be written:

$$\widehat{\boldsymbol{\beta}}^{\text{LASSO}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \mathbf{h}(\mathbf{x}^{(i)})^T \boldsymbol{\beta} \right)^2 + \delta \|\boldsymbol{\beta}\|_1, \quad (1.3.5)$$

with δ a non-negative constant.

The higher δ is, the more zero coefficients there are and the sparser the regression model is.

1.3.4 Forward stagewise regression

Hastie et al. [2001] have introduced the forward stagewise regression. Although different from LASSO, it yields similar results. The procedure can be defined by the following algorithm:

- Initialize with $\mathbf{R} = \mathbf{y}^{\text{obs}}$ and $\beta_i = 0$, $i \in \{1, \dots, p\}$, then repeat until no regressor is correlated with \mathbf{R} :
 - Find $i \in \{1, \dots, p\}$ such that $\mathbf{h}_i(\mathbf{X}^{\text{obs}})$ is the most correlated with \mathbf{R} ,
 - Update $\beta_i = \beta_i + \epsilon_i$, $\epsilon_i = \epsilon \operatorname{sign}(\operatorname{cor}(\mathbf{h}_i(\mathbf{X}^{\text{obs}}), \mathbf{R}))$,
 - Update $\mathbf{R} = \mathbf{R} - \epsilon_i \mathbf{h}_i(\mathbf{X}^{\text{obs}})$,

where, by abuse of notation $\mathbf{h}_i(\mathbf{X}^{\text{obs}}) = (\mathbf{h}_i(\mathbf{x}^{(1)}), \dots, \mathbf{h}_i(\mathbf{x}^{(n)}))$. In practice, ϵ is set to a small value, like $\epsilon = 0.01$. In general, this approach is more reliable than the classical stepwise regression.

1.3.5 Least Angle Regression

Introduced by Efron et al. [2004], Least Angle Regression (LAR) is similar to the forward stagewise regression, given that it selects the regressor $\mathbf{h}_i(\mathbf{X}^{\text{obs}})$ which is the most correlated with the current residual \mathbf{R} . However, the computation of the value of β_i is different. Instead of being slightly modified, the value of β_i is chosen such that the correlation between the new residual $\mathbf{R} - \beta_i \mathbf{h}_i(\mathbf{X}^{\text{obs}})$ and its most correlated regressor $\mathbf{h}_j(\mathbf{X}^{\text{obs}})$ is equal to the correlation between $\mathbf{R} - \beta_i \mathbf{h}_i(\mathbf{X}^{\text{obs}})$ and $\mathbf{h}_i(\mathbf{X}^{\text{obs}})$. This method can also be seen as an intermediate method between forward regression and forward stagewise regression.

1.3.5.1 The algorithm

Least Angle Regression (LAR) is associated with the following algorithm:

1. Initialize with $\mathbf{R} = \mathbf{y}^{\text{obs}}$ and $\beta_i = 0$, $i \in \{1, \dots, p\}$.
2. Find $i \in \{1, \dots, p\}$ such that $\mathbf{h}_i(\mathbf{X}^{\text{obs}})$ is the most correlated with \mathbf{R} .
3. Move β_i from 0 toward its least squares coefficient, until another regressor $\mathbf{h}_j(\mathbf{X}^{\text{obs}})$ has as much correlation with $\mathbf{R} - \beta_i \mathbf{h}_i(\mathbf{X}^{\text{obs}})$ as $\mathbf{h}_i(\mathbf{X}^{\text{obs}})$.
4. Move jointly (β_i, β_j) in the direction defined by their joint least squares coefficient of the current residual on $(\mathbf{h}_i(\mathbf{X}^{\text{obs}}), \mathbf{h}_j(\mathbf{X}^{\text{obs}}))$, until some regressor $\mathbf{h}_k(\mathbf{X}^{\text{obs}})$ is as much correlated with the current residual.
5. Continue until $\min(p, n - 1)$ regressors have been retained.

1.3.5.2 LASSO can be seen as specific case of LAR

Efron et al. [2004] and Hastie et al. [2007] have shown that a slightly modified LAR algorithm can provide the entire paths of the LASSO coefficients as the δ coefficient increases. This modified algorithm is defined as follows:

- Run the LAR algorithm from step 1 to 4,
- If a non-zero coefficient achieves zero, remove the associated regressor from the linear model and recompute the joint least squares direction,
- Continue until $\min(p, n - 1)$ regressors have been retained.

In the same way, a modified LAR algorithm can be used to perform a forward stagewise regression in the case of $\epsilon \rightarrow 0$ [Hastie et al., 2007]. Note that the label LARS generally refers to this modified LAR algorithm (where S refers to Stagewise or LASSO).

1.3.5.3 Hybrid LARS

Introduced by Efron et al. [2004], hybrid LARS is derived from the original LARS (referring to the original LAR or LASSO here). This modified algorithm comprises a LAR step which enables to select the regressors. The next step is the estimation by ordinary least squares of the coefficients associated with the selected regressors.

Hybrid LARS relies on a separation between the choice of the regressors and the estimation of the linear model.

It enables to increase the accuracy of the linear model compared to the original LARS.

Relaxed LASSO [Meinshausen et al., 2007] is an extension of the LARS-based LASSO algorithm. The first step is the same as for hybrid LARS. The ordinary least squares estimation of the coefficients at the second step is replaced by a LASSO estimation with a small penalty. In this approach, for the selected regressors at a given step of the LARS algorithm, one performs LASSO with a small penalty coefficient δ , such that no regressor is eliminated. Hybrid LASSO is a particular case of this algorithm, with $\delta = 0$.

1.3.6 Dantzig selector

The Dantzig selector of Candès and Tao [2007] is based on the resolution of the following optimization problem:

$$\boldsymbol{\beta}^{\text{Dantzig}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{H}^T (y^{\text{obs}} - \mathbf{H}\boldsymbol{\beta})\|_{\infty} \quad \text{subject to } \|\boldsymbol{\beta}\|_1 \leq t, \quad (1.3.6)$$

with $t \in \mathbb{R}^+$

In the same way as LARS, the Dantzig selector sets some coefficients to zero, thus selecting some regressors.

However, Efron et al. [2004] and Meinshausen et al. [2007] have shown that the linear model obtained with LASSO is as accurate as or more accurate than the one obtained with the Dantzig selector.

Note that a DASSO (Dantzig Selector with Sequential Optimization) algorithm has been proposed by James et al. [2008] in order to compute in one step the whole path of the Dantzig selector.

1.3.7 Conclusions

In this section, methods which enable to select the regressors of a linear model have been reviewed. Such approaches are particularly useful when the number of observations n is small compared to the number of possible regressors p of the linear model.

1.4 Gaussian process regression or Kriging

This section is devoted to the surrogate modeling of a computer code by Gaussian Process Regression.

Gaussian process regression is widely used in computer experiments [Sacks et al., 1989; Santner et al., 2003; Rasmussen and Williams, 2006]. In the Gaussian process regression framework, the output y of the code can be seen as a realization of a Gaussian process.

In the remainder of the section, we first outline the multidimensional Gaussian distribution and the definition of a Gaussian process. Then the Gaussian process regression framework for a known covariance function is presented. Finally, the estimation of the hyperparameters of parametric covariance functions is described.

1.4.1 Gaussian processes

1.4.1.1 Multidimensional (multivariate) Gaussian distribution

A random vector $\mathbf{u} = (u_1, \dots, u_n)$, $n \geq 1$, is a Gaussian vector if the following equivalent assumptions are verified:

- for any $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{a}^T \mathbf{u}$ has a Gaussian distribution,
- the characteristic function of \mathbf{u} is of the form $\mathbf{v} \mapsto \exp\left(i\mathbf{v}^T \mathbf{m} - \frac{1}{2}\mathbf{v}^T \mathbf{K} \mathbf{v}\right)$ with \mathbf{m} a n -dimensional vector and \mathbf{K} a $(n \times n)$ -dimensional matrix, which is symmetric and positive definite.

If these assumptions are verified, we have $\mathbf{u} \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$ with $\mathbf{m} = \mathbb{E}[\mathbf{u}]$ and $\mathbf{K} = \operatorname{cov}(\mathbf{u})$.

1.4.1.2 Gaussian processes

A random process associates to any value of \mathbf{x} a random variable $Y(\mathbf{x})$. A random process is a Gaussian process if its finite-dimensional distributions are Gaussian distributions. A Gaussian process Y is characterized by its mean and covariance functions. The mean function is defined by:

$$m(\mathbf{x}) = \mathbb{E}[Y(\mathbf{x})]. \quad (1.4.1)$$

The covariance function is defined by:

$$C(\mathbf{x}, \mathbf{x}') = \text{cov}(Y(\mathbf{x}), Y(\mathbf{x}')), \quad (1.4.2)$$

\mathbf{x}' in \mathbb{X} .

A Gaussian process is said to be stationary if, for all $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ in \mathbb{X} and $\mathbf{h} \in \mathbb{R}^d$ such that $\mathbf{x}^{(1)} + \mathbf{h}, \dots, \mathbf{x}^{(n)} + \mathbf{h}$ are still in \mathbb{X} , the multidimensional distribution of the Gaussian process Y at $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ is the same as the one at $\mathbf{x}^{(1)} + \mathbf{h}, \dots, \mathbf{x}^{(n)} + \mathbf{h}$.

It follows that a covariance function is said to be stationary, if, for all $\mathbf{x}, \mathbf{x}', \mathbf{x} + \mathbf{h}, \mathbf{x}' + \mathbf{h} \in \mathbb{X}$, one has:

$$C(\mathbf{x} + \mathbf{h}, \mathbf{x}' + \mathbf{h}) = C(\mathbf{x}, \mathbf{x}') = C(\mathbf{x} - \mathbf{x}', \mathbf{0}). \quad (1.4.3)$$

Finally, a Gaussian process is stationary if and only if its mean function is constant and its covariance function is stationary.

The next section outlines some classical parametric families of stationary covariance functions and their properties. For a more detailed review of covariance functions, the interested reader may refer to Abrahamsen [1997] and Rasmussen and Williams [2006].

1.4.1.3 Parametric families of stationary covariance functions

Typical parametric families of covariance functions are of the form:

$$C(\mathbf{x}, \mathbf{x}') = \sigma^2 K_{\boldsymbol{\ell}}(\mathbf{x} - \mathbf{x}') \quad (1.4.4)$$

where $K_{\boldsymbol{\ell}}$ is a correlation function parametrized by the vector of correlation lengths $\boldsymbol{\ell} \in (0, +\infty)^d$, and $\sigma^2 \in (0, +\infty)$ is a variance parameter.

The following paragraphs present some classical stationary correlation functions $K_{\boldsymbol{\ell}}$.

The nugget correlation function

The nugget correlation function is defined by:

$$K_{\boldsymbol{\ell}}(\mathbf{x} - \mathbf{x}') = \delta_{\mathbf{x}=\mathbf{x}'}, \quad (1.4.5)$$

where δ denotes the Kronecker delta. Note that this covariance function does not depend on any correlation length.

By construction, the observations of a Gaussian process with a nugget correlation function are not correlated and consequently independent and identically distributed.

Figure 1.1 presents an example of a path of the centered Gaussian process with the nugget correlation function and a unit variance σ^2 . The trajectory is very rough and all the observations are independent of each other.

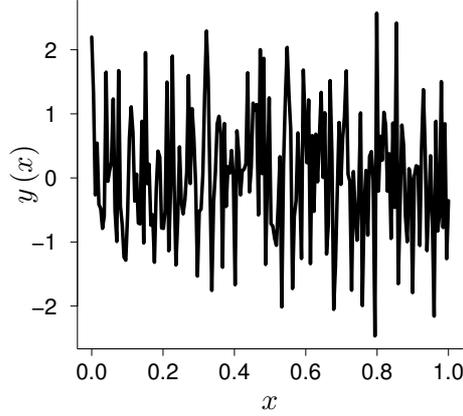


Figure 1.1: An example of a path of the centered Gaussian process with the nugget correlation function and a unit variance σ^2 .

The squared exponential correlation function

The squared exponential (or Gaussian) correlation function is defined by:

$$K_{\ell}(\mathbf{x} - \mathbf{x}') = \exp\left(-d_{\ell}(\mathbf{x} - \mathbf{x}')^2\right), \quad (1.4.6)$$

where $d_{\ell}(\mathbf{x} - \mathbf{x}') = \sqrt{\sum_{i=1}^d \left(\frac{x_i - x'_i}{\ell_i}\right)^2}$. The trajectories of a Gaussian process with a squared exponential correlation function are infinitely differentiable. This covariance function is widely used in Kriging models. However, the assumption of infinite differentiability may be unrealistic [Stein, 1999].

Figure 1.2 presents the squared-exponential correlation function and an example of a path of the centered Gaussian process with a squared-exponential correlation function, a unit variance σ^2 , and the following correlation lengths: $\ell \in \{0.05, 0.1, 0.2\}$. It can be seen that the shorter the correlation length is, the faster the correlation function decreases. Besides, the path varies more if the correlation length is short. Finally, note that the trajectories are very smooth, in agreement with their infinite differentiability.

The Matérn correlation function

The multi-dimensional Matérn kernel can be defined as:

$$K_{\ell}(\mathbf{x} - \mathbf{x}') = \frac{1}{\Gamma(\nu) 2^{\nu-1}} (2\sqrt{\nu}d_{\ell}(\mathbf{x} - \mathbf{x}'))^{\nu} K_{\nu}(2\sqrt{\nu}d_{\ell}(\mathbf{x} - \mathbf{x}')), \quad (1.4.7)$$

with $\Gamma(\cdot)$ the gamma function, K_{ν} a modified Bessel function [Abramowitz and Stegun, 1965] and $\nu \geq \frac{1}{2}$ the smoothness hyperparameter.

Note that as $\nu \rightarrow \infty$, the Matérn kernel tends to the squared exponential correlation function. Besides, when $\nu = k + \frac{1}{2}$, $k \in \mathbb{N}$, the Matérn kernel has a simpler form. In particular, we have:

- if $\nu = \frac{1}{2}$:

$$K_{\ell}(\mathbf{x} - \mathbf{x}') = \exp(-d_{\ell}(\mathbf{x} - \mathbf{x}')), \quad (1.4.8)$$

this kernel is also known as the exponential kernel,

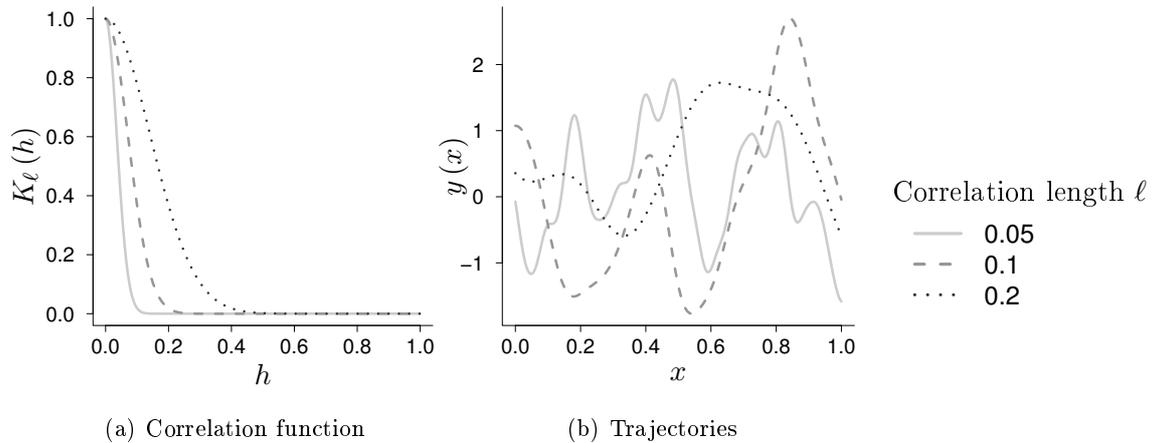


Figure 1.2: On the left figure: plot of the squared-exponential correlation function. On the right plot: an example of a path of the centered Gaussian processes with a squared-exponential correlation function K_ℓ , $\ell \in \{0.05, 0.1, 0.2\}$ and a unit variance.

- if $\nu = \frac{3}{2}$:

$$K_\ell(\mathbf{x} - \mathbf{x}') = \left(1 + \sqrt{3} d_\ell(\mathbf{x} - \mathbf{x}')\right) \exp\left(-\sqrt{3} d_\ell(\mathbf{x} - \mathbf{x}')\right), \quad (1.4.9)$$

- if $\nu = \frac{5}{2}$:

$$K_\ell(\mathbf{x} - \mathbf{x}') = \left(1 + \sqrt{5} d_\ell(\mathbf{x} - \mathbf{x}') + \frac{5}{3} d_\ell(\mathbf{x} - \mathbf{x}')^2\right) \exp\left(-\sqrt{5} d_\ell(\mathbf{x} - \mathbf{x}')\right). \quad (1.4.10)$$

Figure 1.3 presents the exponential correlation function and an example of a path of the centered Gaussian processes with an exponential correlation function, a unit variance σ^2 and the following correlation lengths: $\ell \in \{0.05, 0.1, 0.2\}$. The trajectories are not differentiable.

Figure 1.4 presents the Matérn $\frac{3}{2}$ correlation function and examples of a path of the centered Gaussian processes with a Matérn $\frac{3}{2}$ correlation function, a unit variance σ^2 , and the following correlation lengths: $\ell \in \{0.05, 0.1, 0.2\}$. The trajectories are not very smooth, but smoother than with the exponential correlation function.

Figure 1.5 presents the Matérn $\frac{5}{2}$ correlation function and an example of a path of the centered Gaussian processes with a Matérn $\frac{5}{2}$ correlation function, a unit variance σ^2 , and the following correlation lengths: $\ell \in \{0.05, 0.1, 0.2\}$. The trajectories are relatively smooth.

It can be seen on Figures 1.2 to 1.5 that the shorter the correlation length is, the faster the correlation function decreases. Besides, the path varies more if the correlation length is short.

Figure 1.6 presents the Matérn correlation function and examples of a path of the centered Gaussian processes with a Matérn correlation function, a correlation length equal to 0.5, a unit variance σ^2 , and the following values of the smoothness parameter: $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \infty\}$. It can be seen that the smoothness parameter strongly impacts the form of the correlation function. Besides, the higher ν is, the smoother the paths are.

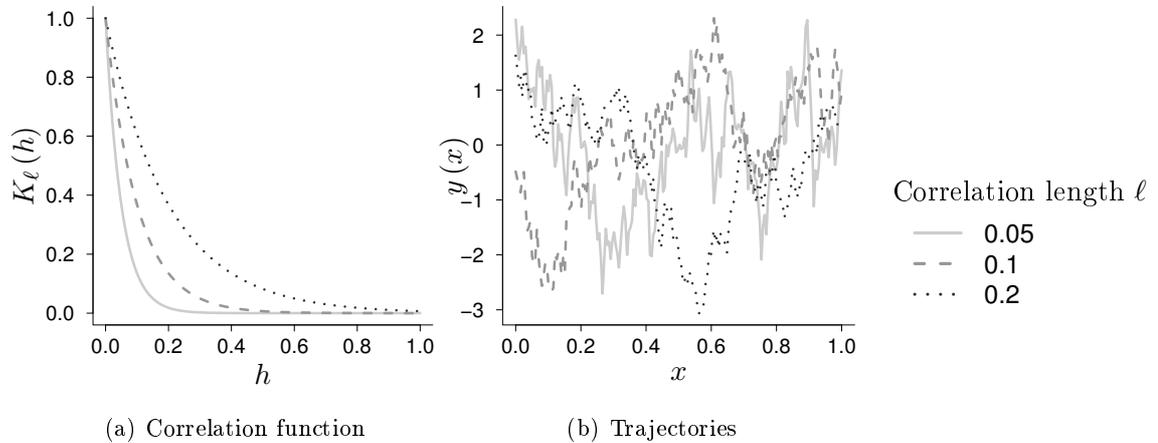


Figure 1.3: On the left figure: plot of the exponential correlation function. On the right plot: an example of paths of the centered Gaussian processes with an exponential correlation function K_ℓ , $\ell \in \{0.05, 0.1, 0.2\}$ and a unit variance.

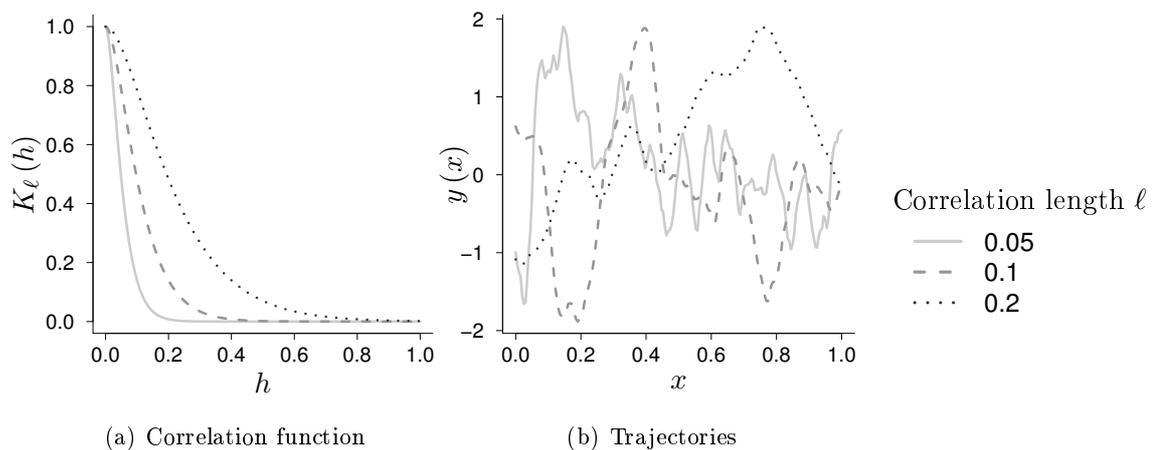


Figure 1.4: On the left figure: plot of the Matérn $\frac{3}{2}$ correlation function. On the right plot: an example of a path of the centered Gaussian processes with a Matérn $\frac{3}{2}$ correlation function K_ℓ , a unit variance and $\ell \in \{0.05, 0.1, 0.2\}$.

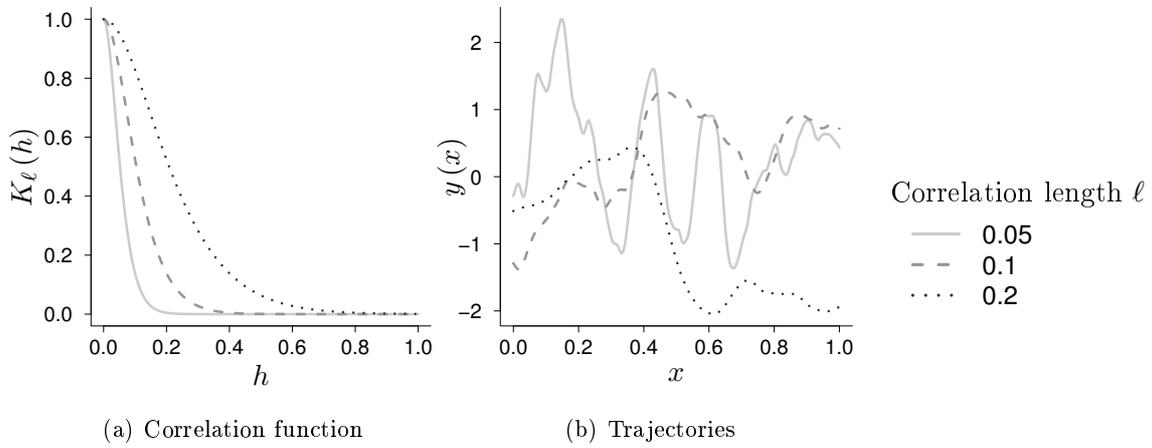


Figure 1.5: On the left figure: plot of the Matérn $\frac{5}{2}$ correlation function. On the right plot: examples of a path of the centered Gaussian processes with a Matérn $\frac{5}{2}$ correlation function K_ℓ , a unit variance and $\ell \in \{0.05, 0.1, 0.2\}$.

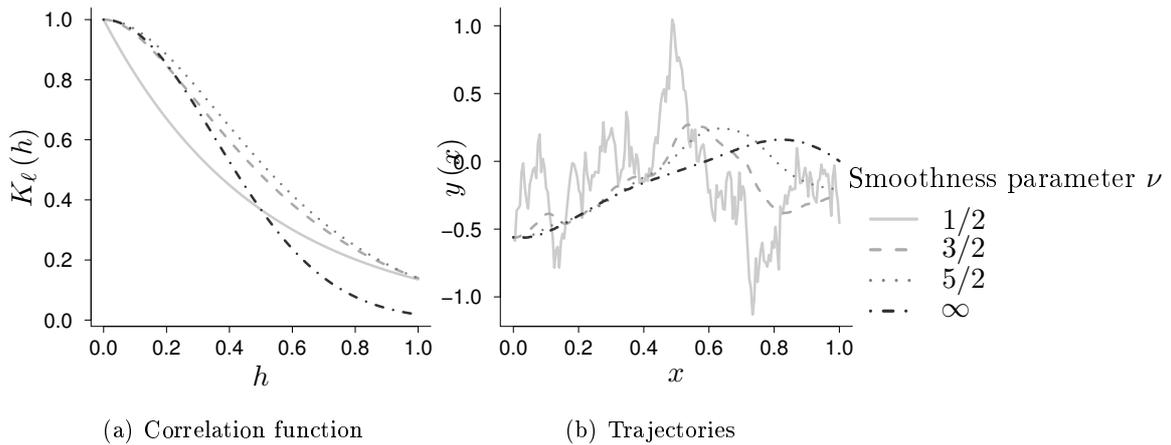


Figure 1.6: On the left figure: plot of the Matérn correlation function for different values of the smoothness parameter ν . On the right plot: an example of a path of the centered Gaussian processes with a Matérn correlation function K_ℓ , the following values of the smoothness parameter $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \infty\}$, a correlation length ℓ equal to 0.5 and a unit variance σ^2 .

The power exponential correlation function

The power exponential correlation kernel is defined by:

$$K_\ell(\mathbf{x} - \mathbf{x}') = \exp\left(-\sum_{i=1}^d \left(\frac{x_i - x'_i}{\ell_i}\right)^p\right), \quad (1.4.11)$$

$p \in (0, 2]$, with the particular case of $p = 2$ corresponding to the squared exponential correlation function.

Finally, note that the multidimensional correlation functions can also be defined as a product of univariate correlation functions:

$$K_\ell(\mathbf{x} - \mathbf{x}') = \prod_{i=1}^d K_{\ell_i}(x_i - x'_i), \quad (1.4.12)$$

where the K_{ℓ_i} may belong to different families of correlation functions.

1.4.1.4 The relationship between the covariance function and the mean square regularity

In this section, we consider a centered Gaussian process Y with covariance function C . Some properties concerning the mean square regularity of a centered Gaussian process and its relationship with the covariance function are reviewed.

A zero-mean Gaussian process Y is mean square continuous if and only if its covariance function is continuous at each pair (\mathbf{x}, \mathbf{x}) , $\mathbf{x} \in \mathbb{X}$. Besides, if a covariance function is continuous at each pair (\mathbf{x}, \mathbf{x}) , $\mathbf{x} \in \mathbb{X}$, then it is continuous on $\mathbb{X} \times \mathbb{X}$ [Bachoc, 2013b].

If one defines the following notation:

$$\text{cov}\left(\frac{\partial Y(\mathbf{x})}{\partial x_i}, \frac{\partial Y(\mathbf{x}')}{\partial x'_i}\right) = \frac{\partial^2 C}{\partial x_i \partial x'_i}(\mathbf{x}, \mathbf{x}'), \quad (1.4.13)$$

the derivative $\frac{\partial}{\partial x_{i_1}} \dots \frac{\partial}{\partial x_{i_k}} Y$, with $\{i_1, \dots, i_k\}$ a subset of $\{1, \dots, d\}$, exists in the mean square sense and is a Gaussian process if the derivative function $\frac{\partial^2}{\partial x_{i_1} \partial x'_{i_1}} \dots \frac{\partial^2}{\partial x_{i_k} \partial x'_{i_k}} C$ exists and is finite.

In the case of a Gaussian process Y with stationary covariance function C , the three following assumptions are a consequence of the previous assumptions:

1. the Fourier transform \widehat{C} of C is such that:

$$\int_{\mathbb{R}} \omega^{2k} \widehat{C}(\omega) d\omega < +\infty,$$

2. the covariance function C of Y is $2k$ times differentiable,
3. Y is k times mean square differentiable.

1.4.2 Ordinary, simple and universal Kriging

The term Kriging [Matheron and Blondel, 1962] refers to the prediction of the value of a random field at unobserved points of this random field. In this work, we assume that the random field is a Gaussian process.

In the framework of Kriging, three cases can be distinguished according to different assumptions on the mean function:

- Simple Kriging corresponds to the case where the mean function is known. Then, thanks to the subtraction of this known mean, the Gaussian process can be assumed to be centered.
- Ordinary Kriging corresponds to the case where the mean function is assumed to be constant and unknown.
- Universal Kriging corresponds to the case where the mean function is unknown and of the form $m(\mathbf{x}) = \mathbf{h}(\mathbf{x})^T \boldsymbol{\beta}$, where $\mathbf{h}(\mathbf{x})$ defines a p -dimensional basis of functions and $\boldsymbol{\beta} \in \mathbb{R}^p$ a vector of unknown coefficients.

Note that, if the covariance function of the Gaussian process is considered as being stationary, the use of a non-stationary mean function (universal Kriging) can make this assumption of stationarity of the covariance function more likely.

In the following paragraphs, we review the predictors obtained by the computation of the conditioned mean and variance of the Gaussian process in the frameworks of simple, ordinary and universal Kriging. At this stage, the covariance function of the Gaussian process is assumed to be known. Besides, we consider a Bayesian framework [Robert, 2007; Santner et al., 2003].

The following notations will be used. The prior distribution of the Gaussian process Y can be denoted by:

$$Y(\cdot) | m, C \sim \text{GP}(m(\cdot), C(\cdot, \cdot)), \quad (1.4.14)$$

and the posterior distribution of the Gaussian process Y by:

$$Y(\cdot) | \mathbf{y}^{\text{obs}}, m, C \sim \text{GP}(m^c(\cdot), C^c(\cdot, \cdot)). \quad (1.4.15)$$

1.4.2.1 Simple Kriging

Simple Kriging corresponds to the case of a Gaussian process with known mean. For the sake of simplicity, this mean is assumed to be set at zero, thanks to the subtraction of the known mean of the Gaussian Process. Thus, one has:

$$m(\mathbf{x}) = 0, \quad (1.4.16)$$

and:

$$Y(\cdot) | C \sim \text{GP}(0, C(\cdot, \cdot)). \quad (1.4.17)$$

In such a case, the conditioned distribution of the Gaussian process is still Gaussian, with conditioned mean and variance which are given by:

$$m^c(\mathbf{x}) = C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{y}^{\text{obs}}, \quad (1.4.18)$$

and

$$C^c(\mathbf{x}, \mathbf{x}') = C(\mathbf{x}, \mathbf{x}') - C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}), \quad (1.4.19)$$

where \mathbf{X}^{obs} is defined by Eq. (1.0.1) and $C(\mathbf{x}, \mathbf{X}^{\text{obs}})$ is a n -dimensional vector and $C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})$ is a $(n \times n)$ -dimensional matrix, so that:

$$(C(\mathbf{x}, \mathbf{X}^{\text{obs}}))_i = C(\mathbf{x}, \mathbf{x}^{(i)}), \quad (1.4.20)$$

and

$$(C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}}))_{ij} = C(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}). \quad (1.4.21)$$

1.4.2.2 Ordinary Kriging

Ordinary Kriging can be regarded as a specific case of Universal Kriging, with constant mean $\beta \in \mathbb{R}$ to be determined:

$$m(\mathbf{x}) = \beta. \quad (1.4.22)$$

Therefore, in the case of Ordinary Kriging, one has:

$$Y(\cdot) | \beta, C \sim \text{GP}(\beta, C(\cdot, \cdot)). \quad (1.4.23)$$

We consider a Bayesian framework and we have no a priori information about β . The prior distribution of β is therefore assumed to be an improper uniform distribution on \mathbb{R} . In such a framework, the conditioned distribution of the Gaussian process is still Gaussian, with the following conditioned mean and variance functions:

$$m^c(\mathbf{x}) = \hat{\beta} + C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} (y(\mathbf{x}) - \hat{\beta}), \quad (1.4.24)$$

and

$$\begin{aligned} C^c(\mathbf{x}, \mathbf{x}') &= C(\mathbf{x}, \mathbf{x}') - C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}) + \\ &\mathbf{u}(\mathbf{x}) \left(\mathbb{1}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbb{1} \right)^{-1} \mathbf{u}(\mathbf{x}'), \end{aligned} \quad (1.4.25)$$

where

$$\mathbf{u}(\mathbf{x}) = \mathbb{1} - \mathbb{1}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}), \quad (1.4.26)$$

$$\hat{\beta} = \left(\mathbb{1}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbb{1} \right)^{-1} \mathbb{1}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{y}^{\text{obs}}, \quad (1.4.27)$$

and:

$$\mathbb{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}. \quad (1.4.28)$$

1.4.2.3 Universal Kriging

In the case of Universal Kriging, the mean function of the Gaussian process is defined as follows:

$$m(\mathbf{x}) = \mathbf{h}(\mathbf{x})^T \boldsymbol{\beta}, \quad (1.4.29)$$

with $\boldsymbol{\beta}$ a vector of unknown parameters.

Therefore, in the case of Universal Kriging, the prior distribution of the Gaussian process is:

$$Y(\cdot) | \mathbf{h}, \boldsymbol{\beta}, C \sim \text{GP}(\mathbf{h}(\cdot)^T \boldsymbol{\beta}, C(\cdot, \cdot)). \quad (1.4.30)$$

If we assume that β follows an improper uniform distribution on \mathbb{R}^p and that the covariance function is known, then the conditional distribution of the Gaussian process is still Gaussian and its conditioned mean and covariance functions can be computed analytically. The conditioned mean and variance of the Gaussian process can be written:

$$m^c(\mathbf{x}) = \mathbf{h}(\mathbf{x})^T \hat{\beta} + C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H}\hat{\beta}), \quad (1.4.31)$$

and

$$C^c(\mathbf{x}, \mathbf{x}') = C(\mathbf{x}, \mathbf{x}') - C(\mathbf{x}, \mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}) + \mathbf{u}(\mathbf{x})^T \left(\mathbf{H}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{H} \right)^{-1} \mathbf{u}(\mathbf{x}'), \quad (1.4.32)$$

where

$$\mathbf{u}(\mathbf{x}) = \mathbf{h}(\mathbf{x}) - \mathbf{H}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}), \quad (1.4.33)$$

and:

$$\hat{\beta} = \left(\mathbf{H}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{y}^{\text{obs}}. \quad (1.4.34)$$

Besides, the posterior distribution of the parameters β is Gaussian with mean $\hat{\beta}$ and covariance:

$$\mathbf{R}_\beta = \left(\mathbf{H}^T C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{H} \right)^{-1}. \quad (1.4.35)$$

Note that the classical linear regression leads to the same results as Universal Kriging with a nugget covariance function. A nugget covariance function is defined by $C(\mathbf{x}, \mathbf{x}') = \sigma^2 \delta_{\mathbf{x}=\mathbf{x}'}$, with δ denoting the Kronecker delta.

1.4.3 Estimation of a parametric covariance function

The previous section has detailed the properties of a Gaussian process and has presented some parametric families of covariance function and the conditioned distribution of the Gaussian process for several assumptions on the mean function of the process and a known covariance function.

In this section, we review some methods of estimation of the hyperparameters of the covariance function, when the covariance function belongs to a known parametric family.

There are two main approaches for the plug-in estimation of the covariance hyperparameters ℓ and σ^2 . The first one is based on the maximization of the likelihood of the observations given the hyperparameters. The second one is based on the minimization of the Leave One Out Mean Square Error for the estimation of ℓ and on the Leave One Out Prediction Variance for the estimation of σ^2 . Alternatively, a full Bayesian approach can be used [Robert, 2007]. But, in such a case the posterior distribution of the Gaussian process is no longer Gaussian.

1.4.3.1 Maximum Likelihood Estimation

By definition of the prior distribution of the Gaussian process modeling the code, one can write:

$$\mathbf{y}^{\text{obs}} \mid \beta, \ell, \sigma^2 \sim \mathcal{N}(\mathbf{H}\beta, \sigma^2 K_\ell(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})), \quad (1.4.36)$$

with K_ℓ such that $C(\mathbf{x}, \mathbf{x}') = \sigma^2 K_\ell(\mathbf{x}, \mathbf{x}')$.

The log-likelihood of the observations can therefore be written as a function of ℓ , σ^2 and β :

$$\mathcal{L}(\beta, \ell, \sigma^2) = -\frac{1}{2} \ln |\sigma^2 \mathbf{R}_\ell| - \frac{1}{2} \frac{1}{\sigma^2} (\mathbf{y}^{\text{obs}} - \mathbf{H}\beta)^T \mathbf{R}_\ell^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H}\beta), \quad (1.4.37)$$

with $\mathbf{R}_\ell = K_\ell(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})$.

The derivatives of $\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\ell}, \sigma^2)$ with respect to $\boldsymbol{\beta}$ and σ^2 are defined as follows:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\beta}}(\boldsymbol{\beta}, \boldsymbol{\ell}, \sigma^2) = \frac{1}{2} \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{R}_\ell^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H}\boldsymbol{\beta}), \quad (1.4.38)$$

$$\frac{\partial \mathcal{L}}{\partial \sigma^2}(\boldsymbol{\beta}, \boldsymbol{\ell}, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2} \frac{1}{\sigma^4} (\mathbf{y}^{\text{obs}} - \mathbf{H}\boldsymbol{\beta})^T \mathbf{R}_\ell^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H}\boldsymbol{\beta}). \quad (1.4.39)$$

From Eqs. (1.4.38) and (1.4.39), it can be inferred that the maximization of the log-likelihood criterion $\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\ell}, \sigma^2)$ with respect to σ^2 and $\boldsymbol{\beta}$ can be solved explicitly. Finally, the Maximum Likelihood estimates of $\boldsymbol{\ell}$, σ^2 and $\boldsymbol{\beta}$ are:

$$\boldsymbol{\beta}_{ML}(\boldsymbol{\ell}) = (\mathbf{H}^T \mathbf{R}_\ell^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}_\ell^{-1} \mathbf{y}^{\text{obs}}, \quad (1.4.40)$$

$$\sigma_{ML}^2(\boldsymbol{\ell}) = \frac{1}{n} (\mathbf{y}^{\text{obs}} - \mathbf{H}\boldsymbol{\beta}_{ML})^T \mathbf{R}_\ell^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H}\boldsymbol{\beta}_{ML}), \quad (1.4.41)$$

$$\boldsymbol{\ell}_{ML} = \underset{\boldsymbol{\ell} \in \mathbb{R}^d}{\text{argmin}} \ln |\sigma_{ML}^2(\boldsymbol{\ell}) \mathbf{R}_\ell| \quad (1.4.42)$$

1.4.3.2 Restricted Maximum Likelihood Estimation

Restricted Maximum Likelihood Estimation (REML) enables to estimate the hyperparameters of the covariance function and the parameters $\boldsymbol{\beta}$ independently. This method is particularly appropriate if the prior distribution of $\boldsymbol{\beta}$ is not a uniform improper distribution. It is based on the left null space of matrix \mathbf{H} . This null space can be associated with a $((n-p) \times n)$ -dimensional matrix \mathbf{W} , such that $\mathbf{W}\mathbf{H} = 0$. If one introduces $\mathbf{w}^{\text{obs}} = \mathbf{W}\mathbf{y}^{\text{obs}}$, one has:

$$\mathbf{w}^{\text{obs}} \sim \mathcal{N}(0, \sigma^2 \mathbf{W} \mathbf{R}_\ell \mathbf{W}^T). \quad (1.4.43)$$

The Restricted Maximum Likelihood can thus be written:

$$\mathcal{L}^{REML}(\boldsymbol{\ell}, \sigma^2) = -\frac{1}{2} \ln |\sigma^2 \mathbf{W} \mathbf{R}_\ell \mathbf{W}^T| - \frac{1}{2} \frac{1}{\sigma^2} (\mathbf{w}^{\text{obs}})^T (\mathbf{W} \mathbf{R}_\ell \mathbf{W}^T)^{-1} \mathbf{w}^{\text{obs}}, \quad (1.4.44)$$

and does not depend on the parameters $\boldsymbol{\beta}$.

1.4.3.3 Cross Validation Estimation

Following Dubrule [1983] and Bachoc [2013b], the correlation length of the covariance function can be estimated by minimizing the Leave One Out Mean Square error. The Leave One Out estimate of the correlation length $\boldsymbol{\ell}$ is given by:

$$\boldsymbol{\ell}_{LOO} = \underset{\boldsymbol{\ell}}{\text{argmin}} MSE_{LOO}, \quad (1.4.45)$$

with

$$MSE_{LOO} = \sum_{i=1}^n \left[(\mathbf{y}^{\text{obs}})_i - m_{-i, \boldsymbol{\ell}}^c(\mathbf{x}^{(i)}) \right]^2, \quad (1.4.46)$$

$m_{-i, \boldsymbol{\ell}}^c(\mathbf{x}^{(i)}) = \mathbb{E} \left[Y(\mathbf{x}^{(i)}) \mid (\mathbf{y}^{\text{obs}})_{-i}, \boldsymbol{\ell} \right]$ and $(\mathbf{y}^{\text{obs}})_{-i}$ denoting all the observations except the i -th observation.

The variance hyperparameter σ^2 can be estimated by setting the value of the Leave-One-Out prediction error to 1. The Leave-One-Out prediction error is defined by:

$$\frac{1}{n} \sum_{i=1}^n \frac{\left((\mathbf{y}^{\text{obs}})_i - m_{-i, \ell_{LOO}}^c(\mathbf{x}^{(i)}) \right)^2}{\sigma^2 K_{-i, \ell_{LOO}}^c(\mathbf{x}^{(i)})}, \quad (1.4.47)$$

with $\sigma^2 K_{-i, \ell_{LOO}}^c(\mathbf{x}^{(i)}) = \mathbb{V} \left[Y(\mathbf{x}^{(i)}) \mid (\mathbf{y}^{\text{obs}})_{-i}, \ell_{LOO}, \sigma^2 \right]$.

Thus, the prediction variance estimate is:

$$\sigma_{LOO}^2 = \frac{1}{n} \sum_{i=1}^n \frac{\left((\mathbf{y}^{\text{obs}})_i - m_{-i, \ell_{LOO}}^c(\mathbf{x}^{(i)}) \right)^2}{K_{-i, \ell_{LOO}}^c(\mathbf{x}^{(i)})}. \quad (1.4.48)$$

Moreover, the two criteria can be evaluated using matrix forms:

$$MSE_{LOO} = \frac{1}{n} (\mathbf{y}^{\text{obs}})^T \tilde{\mathbf{R}}_{\ell}^{-} \text{diag} \left(\tilde{\mathbf{R}}_{\ell}^{-} \right)^{-2} \tilde{\mathbf{R}}_{\ell}^{-} \mathbf{y}^{\text{obs}}, \quad (1.4.49)$$

and:

$$\sigma_{LOO}^2 = \frac{1}{n} (\mathbf{y}^{\text{obs}})^T \tilde{\mathbf{R}}_{\ell}^{-} \text{diag} \left(\tilde{\mathbf{R}}_{\ell}^{-} \right)^{-1} \tilde{\mathbf{R}}_{\ell}^{-} \mathbf{y}^{\text{obs}}, \quad (1.4.50)$$

with $\tilde{\mathbf{R}}_{\ell}^{-} = \mathbf{R}_{\ell}^{-1} - \mathbf{R}_{\ell}^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{R}_{\ell}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}_{\ell}^{-1}$ in the Universal Kriging framework, and $\tilde{\mathbf{R}}_{\ell}^{-} = \mathbf{R}_{\ell}^{-1}$ in the simple Kriging framework.

1.5 Design of experiments

From the previous sections it can be inferred that, by construction, the accuracy of the linear model, in the case of a linear regression, or of the Gaussian process, in the case of Kriging, depends on the choice of the observations. In the following section, we focus on the design of experiments, that is to say the choice of the observations of the code.

1.5.1 Space-filling designs

In this section we focus on space-filling designs. Such designs are adapted to the case of inputs with a uniform distribution on the unit hypercube $[0, 1]^d$. Note that if the inputs have a non-uniform distribution, an isoprobabilistic transformation can be used to make them uniformly distributed over $[0, 1]^d$.

1.5.1.1 LHS designs

Introduced by McKay et al. [1979], Latin Hypercube sampling enables to obtain a sample whose marginals are uniform. If one considers the unit hypercube $[0, 1]^d$, a sample of n points is generated by first dividing each of the d axes of the input domain into n parts. Thus, the unit hypercube is divided into n^d parts and the n observations are drawn uniformly into a selection of n of these small hypercubes. As mentioned above, the n small hypercubes are chosen such that the projection onto each axis leads to exactly n different boxes. Figure 1.7 shows an example of a Latin Hypercube Design.

However, if the projections of a Latin Hypercube design on the marginals are uniformly distributed, the projections of higher dimension are not necessarily uniformly distributed. Two distance-based criteria [Johnson et al., 1990] can be used to characterize the space-filling properties of a design of experiments:

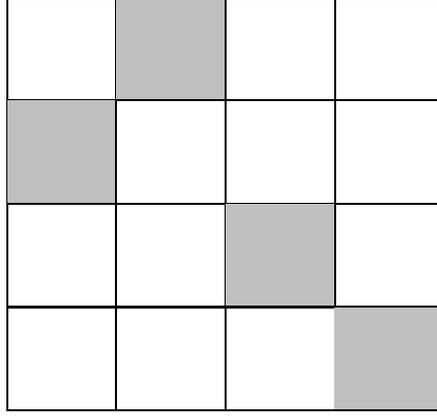


Figure 1.7: An example of Latin Hypercube Design. The observations are drawn in the grey cells.

- the maximin criterion maximizes the Euclidean distance between two points of the design:

$$\mathbf{X}_{\text{maximinLHS}}^{\text{obs}} = \underset{\mathbf{x}^{\text{obs}} \in \mathbb{X}^n}{\operatorname{argmax}} \min_{\substack{i \neq j \\ 1 \leq i, j \leq n}} \left\| (\mathbf{X}^{\text{obs}})_i - (\mathbf{X}^{\text{obs}})_j \right\|, \quad (1.5.1)$$

- the minimax criterion minimizes the distance between any points of \mathbb{X} and the design:

$$\mathbf{X}_{\text{minimaxLHS}}^{\text{obs}} = \underset{\mathbf{x}^{\text{obs}} \in \mathbb{X}^n}{\operatorname{argmin}} \max_{\mathbf{x} \in \mathbb{X}} \max_{1 \leq i \leq n} \left\| \mathbf{x} - (\mathbf{X}^{\text{obs}})_i \right\|. \quad (1.5.2)$$

These criteria can be used in order to sample LHS designs which have good space-filling properties.

1.5.1.2 Quasi-random designs

Low discrepancy sequences like Sobol sequences can also be utilized to ensure good space-filling properties of the design. The notion of discrepancy has been introduced by Niederreiter [1978] and is a measure of the divergence between a set of observations and the uniform distribution. If the definition set is the unit hypercube $[0, 1]^d$, then the discrepancy is defined by:

$$\mathcal{D}(\mathbf{X}^{\text{obs}}) = \sup_{\mathbf{a}, \mathbf{b} \in [0, 1]^d, \mathbf{a} < \mathbf{b}} \left| \frac{\operatorname{card} \left(\left\{ \mathbf{x} \in \mathbf{X}^{\text{obs}} \mid \mathbf{x} \in \prod_{i=1}^d [a_i, b_i] \right\} \right)}{n} - \prod_{i=1}^d (b_i - a_i) \right|, \quad (1.5.3)$$

with $\operatorname{card}(\Omega)$ denoting the number of elements of the finite set Ω .

Low-discrepancy sequences [Niederreiter, 1978] are also known as quasi-random designs. They are defined such that the discrepancy of the sequence tends to zero when the size of the sequence tends to infinity. The low-discrepancy sequences have a smaller discrepancy than a uniform Monte Carlo sample, thus covering better the unit hypercube.

The best-known low-discrepancy sequences are the Van der Corput [Van der Corput, 1935], Halton [Halton, 1964], Sobol [Sobol, 1967], Faure and Hammersley [Hammersley, 1964] sequences.

Space-filling designs can also be defined for the case of a non-hypercube domain (see Perrin

and Cannamela [2017] for example).

If there is no a priori information about the basis of functions in the linear regression case or about the covariance function in the Gaussian process regression case, then space-filling designs are very appropriate to acquire a knowledge of the computer code. Once some information is available, criterion-based designs can be used. The following section details criterion-based designs which are suited for linear regression and Gaussian process regression.

1.5.2 Criterion-based designs

In this section we focus on the optimal designs which can be used when some information about the model is available. The two first sections focus on the criteria which are suited for linear regression and Gaussian process regression. The third section presents the sequential designs, that is to say the enrichment of an initial design (which can be empty) according to a criterion.

1.5.2.1 Designs for linear regression

Elfving [1952] introduced optimal designs for linear regression, with criteria such as D-optimality. Since then, many other criteria, and algorithms of construction of the optimal designs have been proposed [Kiefer and Wolfowitz, 1959; Kiefer, 1961; Fedorov, 1972; Wu and Wynn, 1978; Cook and Nachtsheim, 1980; Fedorov and Hackl, 1997; Molchanov and Zuyev, 2002]. Such designs aim generally at minimizing or maximizing a criterion associated with the variance of the estimation of the regression coefficients β .

According to Eq. (1.4.35), with a nugget covariance of variance σ^2 , the covariance matrix of the posterior distribution of the parameters is given by:

$$\text{cov}(\beta) = \sigma^2 \left(\mathbf{h}(\mathbf{X}^{\text{obs}}) \mathbf{h}(\mathbf{X}^{\text{obs}})^T \right)^{-1}, \quad (1.5.4)$$

where, by abuse of notation, $\mathbf{h}(\mathbf{X}^{\text{obs}})$ is a $(p \times n)$ -dimensional matrix defined by:

$$\mathbf{h}(\mathbf{X}^{\text{obs}}) = \left[\mathbf{h}(\mathbf{x}^{(1)}); \dots; \mathbf{h}(\mathbf{x}^{(n)}) \right]. \quad (1.5.5)$$

Note that the inverse of the covariance matrix of the parameters is also known as the information matrix.

Several criterion-based designs can be used for linear regression:

- The D-optimal criterion aims at maximizing the determinant of the inverse of the covariance matrix:

$$\mathbf{X}_D^{\text{obs}} = \underset{\mathbf{X}^{\text{obs}} \in \mathbb{X}^n}{\text{argmax}} \det \left(\mathbf{h}(\mathbf{X}^{\text{obs}}) \mathbf{h}(\mathbf{X}^{\text{obs}})^T \right), \quad (1.5.6)$$

- The A-optimal criterion aims at minimizing the trace of the covariance matrix:

$$\mathbf{X}_A^{\text{obs}} = \underset{\mathbf{X}^{\text{obs}} \in \mathbb{X}^n}{\text{argmin}} \text{Tr} \left(\left(\mathbf{h}(\mathbf{X}^{\text{obs}}) \mathbf{h}(\mathbf{X}^{\text{obs}})^T \right)^{-1} \right). \quad (1.5.7)$$

1.5.2.2 Designs for Gaussian process regression

In the case of the Gaussian process regression, the design can aim either at improving the estimation of the parameters β of the mean function or at improving the prediction accuracy of the posterior distribution of the Gaussian process.

In the first case, a D-optimal criterion can be used. In the Gaussian process regression framework, this criterion is defined as:

$$\mathbf{X}_D^{\text{obs}} = \operatorname{argmax}_{\mathbf{X}^{\text{obs}} \in \mathbb{X}^n} \det \left(\mathbf{h}(\mathbf{X}^{\text{obs}}) C(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \mathbf{h}(\mathbf{X}^{\text{obs}})^T \right). \quad (1.5.8)$$

If the aim is the improvement of the prediction accuracy, the following criterion, generally referred as I-optimal design, can be used:

$$\mathbf{X}_I^{\text{obs}} = \operatorname{argmin}_{\mathbf{X}^{\text{obs}} \in \mathbb{X}^n} \int_{\mathbb{X}} \mathbb{V}[Y(\mathbf{x}) | \mathbf{y}^{\text{obs}}] d\mu_{\mathbb{X}}(\mathbf{x}). \quad (1.5.9)$$

Note that it can be inferred from Eqs. (1.4.19), (1.4.25) and (1.4.32) that $\mathbb{V}[Y(\mathbf{x}) | \mathbf{y}^{\text{obs}}]$ depends only on \mathbf{X}^{obs} and the covariance function C . By abuse of notation, the previous criterion can be rewritten:

$$\mathbf{X}_I^{\text{obs}} = \operatorname{argmin}_{\mathbf{X}^{\text{obs}} \in \mathbb{X}^n} \int_{\mathbb{X}} \mathbb{V}[Y(\mathbf{x}) | \mathbf{X}^{\text{obs}}] d\mu_{\mathbb{X}}(\mathbf{x}). \quad (1.5.10)$$

The integral $\int_{\mathbb{X}} \mathbb{V}[Y(\mathbf{x}) | \mathbf{X}^{\text{obs}}] d\mu_{\mathbb{X}}(\mathbf{x})$ is defined as the Integrated Mean Square Error (IMSE) [Sacks et al., 1989].

However, the choice of a criterion-based design may pose some difficulties:

- if a discrete search is performed, the number of possible combinations can be very high: $\binom{n}{\mathcal{N}}$, where \mathcal{N} is the number of candidates of the search set.
- in the case of a Gaussian process, the covariance function can be unknown or not precisely known at the beginning.

In those cases, sequential designs can be used. In the case of Gaussian process regression an initial design drawn according to $\mu_{\mathbb{X}}$ can be used for the initial estimation of the covariance function hyperparameters. Then the hyperparameters of the covariance function can be re-estimated at each step of the sequential design.

The stochastic properties of the Gaussian process regression are useful for the definition of sequential designs. Sacks et al. [1989] proposed a sequential design based on the division of the input domain into boxes. The new point is added in the box with the largest contribution to the current IMSE.

Vazquez and Bect [2009] and Bect et al. [2012] proposed a Stepwise Uncertainty Reduction strategy [Geman and Jedynek, 1996] based on a sequential enrichment of the design which is adapted to the estimation of a probability of failure, using a Kriging metamodel and a Bayesian framework.

Such a Stepwise Uncertainty Reduction approach is based on the choice of a new observation point that improves the most a given criterion at the next step.

Bates et al. [1996], then Picheny et al. [2010] proposed a sequential design which is based on the integrated prediction variance (or Integrated Mean Square Error, IMSE) criterion. The associated criterion can be written in the form:

$$\mathbf{x}_{\text{new}} = \operatorname{argmin}_{\mathbf{x}^* \in \mathbb{X}} \int_{\mathbb{X}} \mathbb{V}[Y(\mathbf{x}) | \mathbf{X}^{\text{obs}}, \mathbf{x}^*] d\mu_{\mathbb{X}}(\mathbf{x}), \quad (1.5.11)$$

where, by abuse of notation, $\mathbb{V}[Y(\mathbf{x}) | \mathbf{X}^{\text{obs}}, \mathbf{x}^*] = \mathbb{V}[Y(\mathbf{x}) | \mathbf{y}^{\text{obs}}, y(\mathbf{x}^*)]$. Such a notation can be used, because, for a given covariance function C , the conditioned variance $C^c(\mathbf{x}, \mathbf{x}')$

does not depend on the observations of the output (see Eqs. (1.4.25), (1.4.19) and (1.4.32) for further details).

The above-mentioned design criteria aim at improving the accuracy of the surrogate model, of the posterior distribution of the parameters or of the estimation of a probability of failure. They are all based on the minimization or maximization of a criterion associated with the variance of the estimator of the quantity of interest.

In the next section, we present the Efficient Global Optimization (EGO) algorithm. This is a widely used algorithm which adds to the design a new point which is in the most likely region of a minimum of the function y .

1.5.3 Gaussian processes for pointwise global optimization

Jones et al. [1998] proposed a sequential design aiming at finding the global minimum of an expensive to evaluate function (or computer code). The Efficient Global Optimization algorithm is based on a Gaussian process emulator of the expensive function and takes advantage of the stochastic property of the Gaussian predictor to determine which new point to add. The criterion is based on an Improvement function defined as:

$$I(\mathbf{x}) | \mathbf{y}^{\text{obs}} = \max(\min(\mathbf{y}^{\text{obs}}) - Y(\mathbf{x}) | \mathbf{y}^{\text{obs}}, 0). \quad (1.5.12)$$

The new observation point minimizes the Expected Improvement (EI):

$$\begin{aligned} EI(\mathbf{x}) &= \mathbb{E}[I(\mathbf{x}) | \mathbf{y}^{\text{obs}}] \\ &= (\min(\mathbf{y}^{\text{obs}}) - \mu^c(\mathbf{x})) \Phi\left(\frac{\min(\mathbf{y}^{\text{obs}}) - \mu^c(\mathbf{x})}{\sigma^c(\mathbf{x})}\right) + \sigma^c(\mathbf{x}) \varphi\left(\frac{\min(\mathbf{y}^{\text{obs}}) - \mu^c(\mathbf{x})}{\sigma^c(\mathbf{x})}\right), \end{aligned} \quad (1.5.13)$$

with φ the standard Gaussian probability density function and Φ the standard Gaussian cumulative distribution function. The new observation point \mathbf{x}_{new} is therefore chosen according to the following criterion:

$$\mathbf{x}_{\text{new}} = \underset{\mathbf{x} \in \mathbf{X}}{\operatorname{argmax}} EI(\mathbf{x}). \quad (1.5.14)$$

EGO is a compromise between exploration and exploitation.

1.6 Sensitivity analysis

The sensitivity analysis aims at estimating the importance of the influence of the inputs of a code over the output of the code or over a quantity of interest associated with it. By abuse of notation, this quantity of interest will be denoted by y in the remainder of this section.

The sensitivity analysis methods can be divided into two groups:

- the local sensitivity analysis studies the influence of small variations of the input parameters over a quantity of interest associated with the output of the code,
- the global sensitivity analysis quantifies the influence of the inputs over a quantity of interest associated with the output of the code by considering the variations of the inputs on the whole input domain.

The interested reader can refer to Saltelli et al. [2000] for further details on both groups. In what follows, we will focus on global sensitivity analysis.

Among the methods of global sensitivity analysis, two types of approaches can be distinguished:

- regression-based methods, which are based on the linear regression of the quantity of interest with respect to the inputs. It is worth noticing that such an approach is not adapted to the case of a significantly nonlinear mapping between the inputs and the quantity of interest [Saltelli and Sobol, 1995].
- variance-based methods, which are based on the decomposition of the variance of the quantity of interest with respect to the inputs. This decomposition of variance is also known as ANOVA (Analysis of Variance) [Fisher, 1925]. In particular, the Sobol indices [Sobol, 1993] belong to this category.

In the remainder of the section, we focus on the Sobol indices.

If the variance of the function of interest y is finite and the inputs \mathbf{x} are independent, then the function of interest can be decomposed into first-order effects and interactions [Hoeffding, 1948] :

$$y(\mathbf{x}) = f_0 + \sum_{i=1}^d f_i(x_i) + \sum_{1 \leq i < j \leq d} f_{i,j}(x_i, x_j) + \cdots + f_{1,\dots,d}(\mathbf{x}). \quad (1.6.1)$$

The unique decomposition of y of the form of Eq. (1.6.1) which verifies

$$\text{cov}(f_{i_1,\dots,i_s}(x_{i_1}, \dots, x_{i_s}), f_{j_1,\dots,j_t}(x_{j_1}, \dots, x_{j_t})) = 0, \quad (1.6.2)$$

with $\{i_1, \dots, i_s\} \in \mathbb{N}^s$, $1 \leq i_1 < \dots < i_s \leq d$, $s \in \{1, \dots, d\}$; $\{j_1, \dots, j_t\} \in \mathbb{N}^t$, $1 \leq j_1 < \dots < j_t \leq d$, $t \in \{1, \dots, d\}$ and $\{i_1, \dots, i_s\} \neq \{j_1, \dots, j_t\}$, is defined by [Sobol, 1993] :

$$\begin{aligned} f_0 &= \mathbb{E}[y(\mathbf{x})] \\ f_i(x_i) &= \mathbb{E}[y(\mathbf{x}) | x_i] - f_0 \\ f_{i,j}(x_i, x_j) &= \mathbb{E}[y(\mathbf{x}) | x_i, x_j] - f_i(x_i) - f_j(x_j) - f_0 \\ &\dots \end{aligned} \quad (1.6.3)$$

Given the uncorrelation of the terms of Eq. (1.6.3), the variance of $y(\mathbf{x})$ can thus be decomposed as follows:

$$\mathbb{V}[y(\mathbf{x})] = \sum_{i=1}^d \mathbb{V}[f_i(x_i)] + \sum_{1 \leq i < j \leq d} \mathbb{V}[f_{i,j}(x_i, x_j)] + \cdots + \mathbb{V}[f_{1,\dots,d}(\mathbf{x})], \quad (1.6.4)$$

with the $f_i, f_{i,j} \dots$ defined by Eq. (1.6.3).

The Sobol sensitivity index [Sobol, 1993] corresponding to the subset of input variables $\{x_{i_1}, \dots, x_{i_s}\}$, is defined as:

$$S_{i_1,\dots,i_s} = \frac{\mathbb{V}[\mathbb{E}[y(\mathbf{x}) | x_{i_1}, \dots, x_{i_s}]]}{\mathbb{V}[y(\mathbf{x})]}. \quad (1.6.5)$$

It follows that:

$$1 = \sum_{i=1}^d S_i + \sum_{1 \leq i < j \leq d} S_{i,j} + \cdots + S_{1,\dots,d}. \quad (1.6.6)$$

The first-order Sobol indices are often used to evaluate the individual effect of x_i on y . They are defined as:

$$S_i = \frac{\mathbb{V}[\mathbb{E}[y(\mathbf{x})|x_i]]}{\mathbb{V}[y(\mathbf{x})]}. \quad (1.6.7)$$

Moreover, a total sensitivity index [Homma and Saltelli, 1996] can be defined in order to evaluate the whole contribution of the variable x_i to the variance of the quantity of interest. These total sensitivity indices can be written:

$$T_i = \sum_{\{i_1, \dots, i_s\} \subset \Omega_i} S_{i_1, \dots, i_s}, \quad (1.6.8)$$

where Ω_i denotes the set of all the subsets of $\{1, \dots, d\}$ containing i . These indices can also be written:

$$T_i = 1 - \frac{\mathbb{V}[\mathbb{E}[y(\mathbf{x})|\mathbf{x}_{-i}]]}{\mathbb{V}[y(\mathbf{x})]}, \quad (1.6.9)$$

where \mathbf{x}_{-i} denotes the vector \mathbf{x} except its i -th component .

In practice, the computation of the Sobol indices is performed using Monte Carlo methods [Sobol, 1993]. This computation requires the evaluation of the quantity of interest y at a large number of inputs points. If the computer code associated with this quantity of interest is computationally costly, then the use of a surrogate model of the code becomes necessary [Oakley and O'Hagan, 2004; Le Gratiet, 2013].

Chapter 2

Gaussian process regression of a code with a functional input or output

In this chapter, we review several existing methods for the Gaussian process regression of a computer code with a functional input and/or a functional output. By functional variable, we mean high dimensional vectorial variable. The functional variable is considered to be indexed by the time. The number of indices will be denoted by $N_t \in \mathbb{N}$ in the remainder of this document.

When aiming at performing a Gaussian process regression of a computer code with a functional input, a commonly used approach is to first reduce the dimension of the functional input thanks to a projection technique and then to construct a predictor which is a function of the projection coefficients.

When aiming at performing a Gaussian process regression of the functional output of a computer code with functional output and low dimensional vectorial inputs, two approaches exist. The first one is based on the projection of the output and the independent Gaussian process regression of the projected variables. The second one considers the whole functional output thanks to a tensorized structure of the covariance function of the Gaussian process modeling the code.

This chapter includes therefore two parts. The first one is devoted to the dimension reduction of a functional variable which can be the input or the output of a code. The second one focuses on the Gaussian process regression of the functional output of a code with scalar inputs.

2.1 Dimension reduction of a functional variable

When dealing with functional variables, dimension reduction techniques are often used. In this section, we present some existing methods for the dimension reduction of a functional variable. All the reviewed methods are based on a linear transformation of the functional variable.

The functional variable is denoted by \mathbf{x}_t . Moreover $\mathbf{x}_t \in \mathbb{X}_t \subset \mathbb{R}^{N_t}$, with $N_t \gg 1$, and is associated with the probability measure $\mu_{\mathbb{X}_t}$.

In the considered framework, a set of n observations of the N_t -dimensional vectorial variable \mathbf{x}_t is available. The observations are independently drawn according to $\mu_{\mathbb{X}_t}$ and are centered and gathered in a $(N_t \times n)$ -dimensional matrix $\mathbf{X}_t^{\text{obs}}$:

$$\mathbf{X}_t^{\text{obs}} = \left(\mathbf{x}_t^{(1)} - \overline{\mathbf{x}_t^{\text{obs}}}; \dots; \mathbf{x}_t^{(n)} - \overline{\mathbf{x}_t^{\text{obs}}} \right), \quad (2.1.1)$$

where

$$\bar{\mathbf{x}}_t^{\text{obs}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_t^{(i)}. \quad (2.1.2)$$

Based on this set of observations and for a given dimension of the projection space m , the goal is to find the best m -dimensional set of N_t -dimensional vectors $\{\mathbf{f}_\alpha, \alpha \in \{1, \dots, m\}\}$, and the associated real-valued functions $\mathbf{x}_t \mapsto \beta_\alpha(\mathbf{x}_t)$, which are defined on \mathbb{X}_t . The formalism associated with the dimension reduction of the functional variable \mathbf{x}_t is thus:

$$\mathbf{x}_t \approx \bar{\mathbf{x}}_t^{\text{obs}} + \sum_{\alpha=1}^m \mathbf{f}_\alpha \beta_\alpha(\mathbf{x}_t). \quad (2.1.3)$$

Besides, in the remainder of the section we will consider two types of dimension reduction methods. The first type considers the functional variable only. Such an approach is adapted to the dimension reduction of the functional output of a code, but can also be used for a functional input. The second type reduces the dimension of a functional input \mathbf{x}_t of a code adequately with respect to the output of the code $\mathbf{y}_{\mathbf{x}_t}(\mathbf{x}_t)$.

In this work, we will consider only dimension reductions based on a linear transformation of \mathbf{x}_t and the projection bases are always estimated from the observations.

Note that when considering a code with a functional input, a ridge approximation [Pinkus, 2015; Constantine et al., 2014] can be obtained thanks to the projection of the functional input. Such a ridge approximation can be written in the form:

$$\mathbf{y}_{\mathbf{x}_t}(\mathbf{x}_t) \approx \mathbf{g}_m(\mathbf{B}_m^{\text{obs}}(\mathbf{x}_t - \bar{\mathbf{x}}_t^{\text{obs}})), \quad (2.1.4)$$

where $\mathbf{B}_m^{\text{obs}}$ is a $(m \times N_t)$ -dimensional matrix, \mathbf{g}_m is a function defined on \mathbb{R}^m , whose output has the same dimension as $\mathbf{y}_{\mathbf{x}_t}(\mathbf{x}_t)$, and $\bar{\mathbf{x}}_t^{\text{obs}}$ is defined by Eq. (2.1.2).

In the remainder of the section, we review some methods of dimension reduction of the two types mentioned above:

1. projection of the functional variable which is adapted to the functional variable only,
2. projection of the functional variable which is adapted to a dependent variable.

2.1.1 Methods of dimension reduction based on the functional variable only

When considering only the functional variable and no dependent variable, two types of projection methods can be distinguished. The first type is based on the projection of the functional variable on a basis of *a priori* known functions. The second one, the Principal Components Analysis, relies on the estimation of a projection basis from a set of available observations. These methods can be applied to the case of a functional input or a functional output.

The remainder of this section reviews these two types of approaches.

2.1.1.1 Methods based on the projection on a basis of existing functions

In the case of a basis of existing functions, the vectors \mathbf{f}_α of Eq. (2.1.3) are the discretized versions of functions of time.

The functions can be polynomials, wavelets [Meyer and Salinger, 1995], splines [Hastie et al., 2001], sine and cosine functions...

A set of functions of the basis of size m can be chosen thanks to one of the selection criteria described in Section 1.3. The subset \mathbb{A}_m denotes the indices of the functions which have been kept after the selection procedure.

Moreover, the coefficients $\beta_\alpha(\mathbf{x}_t)$ of Eq. (2.1.3), $\alpha \in \mathbb{A}_m$ can be estimated by solving the following optimization problem:

$$\boldsymbol{\beta}(\mathbf{x}_t) = \underset{\boldsymbol{\beta} \in \mathbb{R}^m}{\operatorname{argmin}} \left\| \mathbf{x}_t - \overline{\mathbf{x}}_t^{\text{obs}} - \mathbf{F}_m \boldsymbol{\beta} \right\|^2, \quad (2.1.5)$$

where \mathbf{F}_m is a $(N_t \times m)$ -dimensional matrix gathering the $\mathbf{f}_\alpha, \alpha \in \mathbb{A}_m$ and $\boldsymbol{\beta}(\mathbf{x}_t)$ is a m -dimensional vector which gathers the $\beta_\alpha(\mathbf{x}_t), \alpha \in \mathbb{A}_m$.

Consequently, $\boldsymbol{\beta}$ is an affine function of \mathbf{x}_t .

The Principal Component Analysis, introduced by Pearson [1901], is a widely used dimension reduction method. It is also known as the Karhunen-Loève expansion [Loève, 1955]. It is based on the eigendecomposition of the covariance matrix of the functional variable. The covariance matrix $\operatorname{cov}(\mathbf{x}_t)$ can be estimated from the set of observations of the functional variable $\mathbf{X}_t^{\text{obs}}$, where $\mathbf{X}_t^{\text{obs}}$ is defined by Eq. (2.1.1). This estimate of the covariance matrix is thus given by:

$$\mathbf{R}_{\mathbf{x}_t}^{\text{obs}} = \frac{1}{n-1} \mathbf{X}_t^{\text{obs}} (\mathbf{X}_t^{\text{obs}})^T. \quad (2.1.6)$$

The projection basis is then defined by the eigenvectors of the covariance matrix $\mathbf{R}_{\mathbf{x}_t}^{\text{obs}}$. In other words, if the eigendecomposition of $\mathbf{R}_{\mathbf{x}_t}^{\text{obs}}$ is denoted by:

$$\mathbf{R}_{\mathbf{x}_t}^{\text{obs}} = \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^T, \quad (2.1.7)$$

where the diagonal of $\boldsymbol{\Lambda}$ gathers the positive decreasing eigenvalues of $\mathbf{R}_{\mathbf{x}_t}^{\text{obs}}$, the m first projected variables are $\mathbf{V}_m^T \mathbf{x}_t$, with \mathbf{V}_m gathering the m first columns of \mathbf{V} .

Note that the accuracy of the approximation $\mathbf{R}_{\mathbf{x}_t}^{\text{obs}}$ of the covariance matrix $\operatorname{cov}(\mathbf{x}_t)$ and thus the one of the projection basis, depend on the available observations of the functional variable.

It is also worth noticing that the Principal Component Analysis can also be used in combination with a ridge approximation. In such a case, the matrix $\mathbf{B}_m^{\text{obs}}$ of Eq. (2.1.4) is equal to \mathbf{V}_m^T .

2.1.2 Methods of dimension reduction of a functional variable which are adapted to a dependent variable

In this section, we focus on linear transformations aiming at reducing the dimension of the functional input of a code, such that the projected variable is adapted to the output of the code.

The two parts of this section present two methods of dimension reduction: the first one is based on Partial Least Squares [Wold, 1966] and the second one is based on the Active Subspaces method [Russi, 2010].

2.1.2.1 Partial Least Squares

Introduced by Wold [1966], Partial Least Squares aim at reducing the dimension of a functional variable \mathbf{x}_t by taking into account a dependent variable which can be a scalar variable $y_{\mathbf{x}_t}$ or a functional variable $\mathbf{y}_{\mathbf{x}_t}$. In our framework, this dependent variable is the output of the code, whereas \mathbf{x}_t is the input. The projection basis is determined from the covariance matrix between the functional variable \mathbf{x}_t and the dependent variable. In this way, the functional input can be projected on a basis which is adapted to the output.

If a set of observations of the output of the code is available, and is denoted by:

$$\mathbf{Y}_{\mathbf{x}_t}^{\text{obs}} = \left(\mathbf{y}_{\mathbf{x}_t} \left(\mathbf{x}_t^{(1)} \right); \dots; \mathbf{y}_{\mathbf{x}_t} \left(\mathbf{x}_t^{(n)} \right) \right), \quad (2.1.8)$$

where $\mathbf{Y}_{\mathbf{x}_t}^{\text{obs}}$ is a $(N_y \times n)$ -dimensional matrix, N_y is the dimension of the output of $\mathbf{y}_{\mathbf{x}_t}$, then the covariance matrix $\text{cov}(\mathbf{x}_t, \mathbf{y}_{\mathbf{x}_t})$ can be approximated by:

$$\mathbf{R}_{\mathbf{x}_t, \mathbf{y}_{\mathbf{x}_t}}^{\text{obs}} = \frac{1}{n-1} \mathbf{X}_t^{\text{obs}} \left(\mathbf{Y}_{\mathbf{x}_t}^{\text{obs}} - \overline{\mathbf{y}_{\mathbf{x}_t}}^{\text{obs}} \right)^T, \quad (2.1.9)$$

where

$$\overline{\mathbf{y}_{\mathbf{x}_t}}^{\text{obs}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_{\mathbf{x}_t} \left(\mathbf{x}_t^{(i)} \right). \quad (2.1.10)$$

Following Höskuldsson [1988], if the singular value decomposition of the covariance matrix between the functional variable and the dependent variable is denoted by:

$$\mathbf{R}_{\mathbf{x}_t, \mathbf{y}_{\mathbf{x}_t}}^{\text{obs}} = \mathbf{U} \mathbf{D} \mathbf{V}^T, \quad (2.1.11)$$

where the diagonal of \mathbf{D} gathers the positive singular values in decreasing order, then the m first projected variables of the functional input which are adapted to the output of the code are given by $\mathbf{U}_m^T \mathbf{x}_t$, with \mathbf{U}_m gathering the m first columns of \mathbf{U} .

If we refer to the ridge approximation of Eq. (2.1.4), then, in the case of Partial Least Squares, $\mathbf{B}_m^{\text{obs}} = \mathbf{U}_m^T$.

Note that the accuracy of the estimation of the covariance matrix $\text{cov}(\mathbf{x}_t, \mathbf{y}_{\mathbf{x}_t})$ and thus the one of the projection basis depend on the number of observations of the functional variable and its dependent variable.

Finally, Nanty et al. [2017] have studied ridge approximations based on the conditioned mean of a Gaussian process (also known as Kriging of Gaussian process regression, see section 1.4 for further details) indexed by the projection of the functional input. They have compared the prediction accuracy with a projection based on Principal Components Analysis or on Partial Least Squares. It is shown that, in many cases, the prediction accuracy of the ridge approximation is better with a dimension reduction based on Partial Least Squares than on Principal Components Analysis.

2.1.2.2 Active Subspaces

Introduced by Russi [2010], the Active Subspace refers to the projection of the functional input on an "Active Subspace", estimated from the observations of the derivatives of the output of the code with respect to the functional input \mathbf{x}_t . Using the formalism introduced by Constantine et al. [2014] for $N_y = 1$, if the set of n observations of the derivatives is denoted by:

$$\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} = \left(\nabla y_{\mathbf{x}_t} \left(\mathbf{x}_t^{(1)} \right); \dots; \nabla y_{\mathbf{x}_t} \left(\mathbf{x}_t^{(n)} \right) \right), \quad (2.1.12)$$

where $\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}}$ is a $(N_t \times n)$ -dimensional matrix, then the projection basis is given by the eigenvectors of the $(N_t \times N_t)$ -dimensional matrix $\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} \left(\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} \right)^T$. In other words, if one denotes by:

$$\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} \left(\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} \right)^T = \mathbf{W} \boldsymbol{\lambda} \mathbf{W}^T \quad (2.1.13)$$

the eigendecomposition of the matrix $\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}} (\nabla \mathbf{y}_{\mathbf{x}_t}^{\text{obs}})^T$, where the diagonal of $\boldsymbol{\lambda}$ gathers the eigenvalues in decreasing order, then the m -dimensional vector of projection coefficients of the functional input \mathbf{x}_t is given by $\mathbf{W}_m^T \mathbf{x}_t$ where \mathbf{W}_m gathers the m first columns of the matrix \mathbf{W} .

If a ridge approximation is performed, the projection matrix $\mathbf{B}_m^{\text{obs}}$, defined by Eq. (2.1.4) is such that $\mathbf{B}_m^{\text{obs}} = \mathbf{W}_m$.

Zahm et al. compare the ridge approximation of $\mathbf{y}_{\mathbf{x}_t}$ for the case of $N_y > 1$ with a projection of the functional input based either on Principal Components Analysis (also called Karhunen-Loève expansion) or on Active Subspaces.

The Active Subspace, given by the projector P_m , is computed from the following matrix:

$$\frac{1}{n} \sum_{i=1}^n \nabla \mathbf{y}_{\mathbf{x}_t} (\mathbf{x}_t^{(i)})^T \nabla \mathbf{y}_{\mathbf{x}_t} (\mathbf{x}_t^{(i)}) \quad (2.1.14)$$

with $\nabla \mathbf{y}_{\mathbf{x}_t} (\mathbf{x}_t^{(i)})$ the $(N_t \times N_y)$ -dimensional matrix of the derivatives at $\mathbf{x}_t^{(i)}$.

The studied ridge approximation of $\mathbf{y}_{\mathbf{x}_t} (\mathbf{x}_t)$ is of the form $\mathbb{E} [\mathbf{y}_{\mathbf{x}_t} (P_m \mathbf{x}_t + P_m^c \mathbf{X}_t) | \mathbf{x}_t]$, where P_m is a projector from \mathbb{R}^{N_t} to \mathbb{R}^m , P_m^c its complement, and \mathbf{X}_t is N_t -dimensional vector with probability measure $\mu_{\mathbb{X}_t}$. The authors conclude that Active Subspaces can yield more effective dimension reduction for the ridge approximation than Principal Components Analysis. They also observe that, if there is no low dimensional structure in the input-output map, then a dimension reduction based on the covariance of the input only (PCA) is more efficient.

This section has been devoted to the dimension reduction of a functional variable, which can be the input or the output of a computer code. In the next section, we focus on the Gaussian process regression of the functional output of a computer code. The notations used will be similar to those of Chapter 1.

2.2 Gaussian process prediction of a computer code with a functional output

In this section, we consider a computer code with low dimensional vectorial inputs and a functional output, that is to say, of the form $\mathbf{x} \mapsto \mathbf{y}(\mathbf{x})$, $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^{N_t}$, $N_t \gg 1$. Moreover, $\mu_{\mathbb{X}}$ is a probability measure associated with \mathbf{x} .

The following sections detail the state of the art for the Gaussian process regression of \mathbf{y} from a set of n observations of the input and the output of the code. These observations are denoted by:

$$\mathbf{X}^{\text{obs}} = \begin{pmatrix} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{x}^{(n)} \end{pmatrix}, \quad (2.2.1)$$

and

$$\mathbf{Y}^{\text{obs}} = \left(\mathbf{y}^{(1)} = \mathbf{y}(\mathbf{x}^{(1)}); \dots; \mathbf{y}^{(n)} = \mathbf{y}(\mathbf{x}^{(n)}) \right), \quad (2.2.2)$$

where \mathbf{X}^{obs} is a $(n \times d)$ -dimensional matrix and \mathbf{Y}^{obs} is a $(N_t \times n)$ -dimensional matrix.

The first subsection of this section focuses on the Gaussian process prediction of a functional output thanks to the projection of this output on a basis. The second subsection is devoted to the Gaussian process regression of the whole functional output of the code.

2.2.1 Projection of the functional output on a basis

Bayarri et al. [2007] proposed to use a wavelet decomposition as a basis representation of the functional output. A thresholding procedure is performed in order to reduce the size of the set of the projection functions while obtaining an accurate projection. Then independent Gaussian predictors of each of the coefficients of the retained projection functions are constructed.

Higdon et al. [2008] proposed to build a Gaussian process emulator of the functional output of a code through a Principal Component Analysis of the functional output. First, a Principal Component Analysis of the functional output is performed. A number m of the projected variables is chosen such that these m components represent 99% of the total variance of the output. If $\overline{\mathbf{y}}^{\text{obs}}$ is the empirical mean of the observations of the output of the code:

$$\overline{\mathbf{y}}^{\text{obs}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}^{(i)}, \quad (2.2.3)$$

where the $\mathbf{y}^{(i)}$ are defined in Eq. (2.2.2), then the functional output of the code can be approximated by:

$$\mathbf{y}(\mathbf{x}) \approx \overline{\mathbf{y}}^{\text{obs}} + \mathbf{V}_m \boldsymbol{\omega}(\mathbf{x}), \quad (2.2.4)$$

with \mathbf{V}_m a $(N_t \times m)$ -dimensional matrix, whose columns are the m first eigenvectors of the empirical covariance matrix $\text{cov}(\mathbf{Y}^{\text{obs}})$ and $\boldsymbol{\omega}$ a m -dimensional function giving the projection coefficients.

The observations of the function giving the projection coefficients are given by:

$$\boldsymbol{\omega}^{\text{obs}} = \mathbf{V}_m^T \left(\mathbf{Y}^{\text{obs}} - \overline{\mathbf{y}}^{\text{obs}} \right), \quad (2.2.5)$$

and $\boldsymbol{\omega}^{\text{obs}}$ is a $(m \times n)$ -dimensional matrix.

Note that, by construction, $\boldsymbol{\omega}$ is expected to be a zero-mean m -dimensional vector.

The components of the function $\boldsymbol{\omega}$ are treated as being independent and a predictor of each component is constructed using the simple Kriging framework:

$$\omega_i(\cdot) \sim \text{GP}(0, C_{\omega_i}(\cdot, \cdot)), \quad (2.2.6)$$

with C_{ω_i} a covariance function, and $1 \leq i \leq m$.

The posterior predictor of the i -th component of function $\boldsymbol{\omega}$ is given by:

$$\omega_i(\cdot) | \boldsymbol{\omega}_i^{\text{obs}} \sim \text{GP}(\mu_{\omega_i}^c(\cdot), C_{\omega_i}^c(\cdot, \cdot)), \quad (2.2.7)$$

where $\boldsymbol{\omega}_i^{\text{obs}}$ corresponds to the i -th line of $\boldsymbol{\omega}^{\text{obs}}$, and:

$$\mu_{\omega_i}^c(\mathbf{x}) = C_{\omega_i}^c(\mathbf{x}, \mathbf{x}^{\text{obs}}) C_{\omega_i}^c(\mathbf{x}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} \boldsymbol{\omega}_i^{\text{obs}}, \quad (2.2.8)$$

and:

$$C_{\omega_i}^c(\mathbf{x}, \mathbf{x}') = C_{\omega_i}(\mathbf{x}, \mathbf{x}') - C_{\omega_i}(\mathbf{x}, \mathbf{X}^{\text{obs}}) C_{\omega_i}(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})^{-1} C_{\omega_i}(\mathbf{X}^{\text{obs}}, \mathbf{x}'), \quad (2.2.9)$$

where $C_{\omega_i}(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})$ is a $(n \times n)$ -dimensional matrix such that

$$C_{\omega_i}(\mathbf{X}^{\text{obs}}, \mathbf{X}^{\text{obs}})_{kl} = C_{\omega_i}(\mathbf{x}^{(k)}, \mathbf{x}^{(l)}), \quad (2.2.10)$$

and $C_{\omega_i}(\mathbf{x}, \mathbf{X}^{\text{obs}})$ is a n -dimensional vector such that:

$$C_{\omega_i}(\mathbf{x}, \mathbf{X}^{\text{obs}})_k = C_{\omega_i}(\mathbf{x}, \mathbf{x}^{(k)}). \quad (2.2.11)$$

The multivariate predictor of $\boldsymbol{\omega}$ is therefore defined by:

$$\boldsymbol{\omega}(\cdot) | \boldsymbol{\omega}^{\text{obs}} \sim \text{GP}(\boldsymbol{\mu}_{\boldsymbol{\omega}}^c(\cdot), \mathbf{C}_{\boldsymbol{\omega}}^c(\cdot, \cdot)), \quad (2.2.12)$$

where:

$$(\boldsymbol{\mu}_{\boldsymbol{\omega}}^c(\mathbf{x}))_i = \mu_{\omega_i}^c(\mathbf{x}) \quad (2.2.13)$$

and:

$$(\mathbf{C}_{\boldsymbol{\omega}}^c(\mathbf{x}, \mathbf{x}'))_{ij} = C_{\omega_i}^c(\mathbf{x}, \mathbf{x}') \delta_{i=j}. \quad (2.2.14)$$

Finally, a predictor of \mathbf{y} is given by:

$$\mathbf{y}(\cdot) | \mathbf{Y}^{\text{obs}} \sim \text{GP}\left(\overline{\mathbf{y}}^{\text{obs}} + \mathbf{V}_m \boldsymbol{\mu}_{\boldsymbol{\omega}}^c(\cdot), \mathbf{V}_m \mathbf{C}_{\boldsymbol{\omega}}^c(\cdot, \cdot) \mathbf{V}_m^T\right). \quad (2.2.15)$$

Perrin [2018] mentions that if the projection basis is estimated from a small set of observations, its estimation may be not very accurate. The accuracy of the prediction of the functional output using the method described above can thus suffer from this lack of accuracy of the projection basis.

2.2.2 Gaussian process regression of the whole functional output

Another possible approach for the Gaussian process regression of a functional output is to choose an appropriate structure of the covariance function of the Gaussian process. Such an approach enables to emulate the whole functional output of a code.

Williams et al. [2006] proposed to treat the index of the functional input as one of the inputs of the model. The covariance function of the Gaussian process depends on the inputs of the code and on the index of the functional output. The output can therefore be treated as a univariate output, indexed by an index input. A power exponential covariance function is used, such that the covariance function has a tensorized structure between the index (time) and the other inputs.

Rougier [2008] and Conti et al. [2009] have used a tensorized structure for the mean and covariance functions of the process. In this framework, the functional output of the code \mathbf{y} can be seen as a Gaussian process \mathbf{Y} with the following properties:

$$\mathbf{Y}(\cdot) | \mathbf{M}, \mathbf{R}_t, C \sim \text{GP}(\mathbf{M}\mathbf{h}(\cdot), \mathbf{R}_t \otimes C(\cdot, \cdot)), \quad (2.2.16)$$

with \mathbf{M} a $(N_t \times p)$ -dimensional matrix, \mathbf{h} a vector of p basis functions, \mathbf{R}_t a $(N_t \times N_t)$ -dimensional covariance matrix and C a covariance function, and \otimes denoting the Kronecker product.

In this framework, if \mathbf{M} has a uninformative prior distribution given by the uniform distribution on the space of the real-valued $(N_t \times p)$ -dimensional matrices, then the distribution of \mathbf{M} given the observations is Gaussian, with the following mean:

$$\begin{aligned} \widehat{\mathbf{M}} &= \mathbb{E}[\mathbf{M} | \mathbf{y}^{\text{obs}}, \mathbf{R}_t, C] \\ &= \mathbb{E}[\mathbf{M} | \mathbf{y}^{\text{obs}}, C] \\ &= \mathbf{y}^{\text{obs}} (\mathbf{R}^{\text{obs}})^{-1} (\mathbf{H}^{\text{obs}})^T \left(\mathbf{H}^{\text{obs}} (\mathbf{R}^{\text{obs}})^{-1} (\mathbf{H}^{\text{obs}})^T \right)^{-1}, \end{aligned} \quad (2.2.17)$$

where \mathbf{R}^{obs} is a $(n \times n)$ -dimensional matrix such that:

$$(\mathbf{R}^{\text{obs}})_{kl} = C(\mathbf{x}^{(k)}, \mathbf{x}^{(l)}), \quad (2.2.18)$$

2.2. GAUSSIAN PROCESS PREDICTION OF A COMPUTER CODE WITH A FUNCTIONAL OUTPUT

and \mathbf{H}^{obs} is a $(p \times n)$ -dimensional matrix whose j -th column is given by $\mathbf{h}(\mathbf{x}^{(j)})$. From Eq. (2.2.16), it can be inferred that:

$$\mathbf{Y}^{\text{obs}} | \mathbf{M}, \mathbf{R}_t, C \sim \mathcal{N}(\mathbf{M}\mathbf{H}^{\text{obs}}, \mathbf{R}_t \otimes \mathbf{R}^{\text{obs}}). \quad (2.2.19)$$

Therefore, the matrix \mathbf{R}_t can be estimated by maximizing the likelihood of the observations, as proposed in Perrin [2018]:

$$\hat{\mathbf{R}}_t = \frac{1}{n} \left(\mathbf{Y}^{\text{obs}} - \widehat{\mathbf{M}}\mathbf{H}^{\text{obs}} \right) (\mathbf{R}^{\text{obs}})^{-1} \left(\mathbf{Y}^{\text{obs}} - \widehat{\mathbf{M}}\mathbf{H}^{\text{obs}} \right)^T. \quad (2.2.20)$$

Finally, in the Universal Kriging framework, with an improper uniform prior for \mathbf{M} , the conditioned distribution of \mathbf{Y} is given by:

$$\mathbf{Y}^c(\cdot) := \mathbf{Y}(\cdot) | \mathbf{Y}^{\text{obs}}, C \sim \text{GP} \left(\boldsymbol{\mu}^c(\cdot), \hat{\mathbf{R}}_t \otimes C^c(\cdot, \cdot) \right), \quad (2.2.21)$$

where:

$$\begin{aligned} \boldsymbol{\mu}^c(\mathbf{x}) &= \widehat{\mathbf{M}}\mathbf{h}(\mathbf{x}) + \left[\mathbf{Y}^{\text{obs}} - \widehat{\mathbf{M}}\mathbf{H}^{\text{obs}} \right] (\mathbf{R}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}), \\ C^c(\mathbf{x}, \mathbf{x}') &= C(\mathbf{x}, \mathbf{x}') - C(\mathbf{x}, \mathbf{X}^{\text{obs}}) (\mathbf{R}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}') \\ &\quad + \mathbf{u}(\mathbf{x})^T \left(\mathbf{H}^{\text{obs}} (\mathbf{R}^{\text{obs}})^{-1} (\mathbf{H}^{\text{obs}})^T \right)^{-1} \mathbf{u}(\mathbf{x}'), \\ \mathbf{u}(\mathbf{x}) &= \mathbf{h}(\mathbf{x}) - \mathbf{H}^{\text{obs}} (\mathbf{R}^{\text{obs}})^{-1} C(\mathbf{X}^{\text{obs}}, \mathbf{x}), \end{aligned} \quad (2.2.22)$$

and $C(\mathbf{x}, \mathbf{X}^{\text{obs}})$ is a n -dimensional vector, such that:

$$C(\mathbf{x}, \mathbf{X}^{\text{obs}})_k = C(\mathbf{x}, \mathbf{x}^{(k)}). \quad (2.2.23)$$

Part II

Contributions

Chapter 3

Nested polynomial trends for the improvement of Gaussian predictors

In this chapter, we focus on the case of two nested codes with scalar outputs. Moreover, there are no observations of the intermediary variable. We therefore consider the following system:

$$\begin{array}{ccc} & \mathbf{x}_2 & \\ & \searrow & \\ \mathbf{x}_1 & \rightarrow & y_1(\mathbf{x}_1) & \nearrow & y_{\text{nest}}(\mathbf{x}_{\text{nest}}) := y_2(y_1(\mathbf{x}_1), \mathbf{x}_2), & (3.0.1) \end{array}$$

where \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_{nest} are low dimensional vectors and y_1 , y_2 and y_{nest} are scalars. This system becomes therefore:

$$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \rightarrow y(\mathbf{x}) := y_2(y_1(\mathbf{x}_1), \mathbf{x}_2). \quad (3.0.2)$$

The work presented in this chapter has been published in [Perrin et al., 2017]. The framework of Gaussian process regression is considered (see Chapter 1 for further details). An innovative parametrization of the mean function of the Gaussian process, based on the composition of two polynomials, is proposed.

3.1 Introduction

The numerical cost of many codes to simulate complex physical systems is very high. In order to perform sensitivity analyses, uncertainty quantification or reliability studies, these computer models have therefore to be replaced by surrogate models, that is to say by fast and inexpensive mathematical functions. Within the computational science community, when the maximal available information is a finite set of code evaluations, the most widely used surrogate models are the generalized polynomial chaos expansion (PCE) [Ghanem and Spanos, 1990, 2003; Soize and Ghanem, 2004; Das et al., 2009; Le Maître and Knio, 2010; Arnst et al., 2010; Perrin et al., 2012] and the Gaussian process regression (GPR), or Kriging (see Sacks et al. [1989]; Oakley and O'Hagan [2002]; Rasmussen and Williams [2006]).

On the one hand, the main idea of PCE is to expand the code output, which is denoted by y in the following, onto an appropriate basis made of orthonormal multivariate polynomials, which are related to the distribution of the code input variables. As the number of unknown expansion coefficients usually grows exponentially with the number of input parameters, the relevance of these approaches strongly depends on their ability to select the most relevant basis functions. To this end, several penalization techniques, such as the ℓ_1 -minimization [Tibshirani, 1989; Jakeman et al., 2015] and the least Angle Regression (LAR) methods [Hastie et al., 2002; Efron et al., 2004; Blatman and Sudret, 2011], have been introduced to select

polynomial basis sets that lead to more accurate PCE than would have been obtained if the basis is *a priori* fixed. Taking advantage of the tensor-product structure of the multivariate polynomial basis, separated representations, such as low-rank approximations [Nouy, 2010; Konakli and Sudret, 2016], have alternatively been proposed to develop surrogate models with polynomial functions in highly-compressed formats.

On the other hand, the GPR is based on the assumption that the code output is a particular realization of a Gaussian stochastic process, Y . This hypothesis, which was first introduced in time series analysis [Parzen, 1962] and in optimization [Kushner, 1964], is widely used as it allows dealing with the conditional probability and expectation, while leading to very interesting results in terms of computer code prediction. Hence, contrary to the PCE, the GPR is not associated with an *a priori* projection basis, but requires the introduction of the mean and the covariance functions of Y . In practice, we observe that the role of the mean function of Y on the prediction decreases when the number of code evaluations increases. This explains that in applications where many code evaluations are available, good GPR-based surrogate models can be obtained using constant or linear trends for the mean function. On the contrary, when the number of code evaluations is small compared to the complexity of y , it can be very useful to optimize it. In that case, searching the mean function of Y as a well-chosen sum of polynomial functions can indeed strongly improve the relevance of the associated GPR. In particular, the authors refer to [Joseph et al., 2008] and [Kersaudy et al., 2015] for an illustration of the interest of using variable selection techniques to optimize this polynomial representation of the mean function of Y .

Following these works, the idea of this part is to propose an alternative parametrization of the mean function of Y , which is particularly adapted to the case when the number of code evaluations is small compared to the complexity of y . Instead of searching sparse polynomial approximations, we look for high dimensional polynomial approximations that are characterized by a small number of parameters. In other words, if we want to model a complex code response with a very limited number of code evaluations, we believe that it can be more efficient to use complex but approximated models than simple but fully optimized models. We thus propose to consider the composition of two polynomials for the mean function of Y . Indeed, the composition of two polynomial functions is still a polynomial function, but of much higher order. In particular, such a formalism can be used to model separately a transformation of each code input and the dependence structure between them.

The main difficulty concerning this specific representation is the identification of the parameters of the two combined polynomials. Indeed, by composing two polynomial functions that are linear with respect to their parameters, we get a strongly non-linear representation, which is likely to be very sensitive to small changes in the parameters' values. In addition, distinct values for these parameters can lead to the same nested representation, which does not help for the identification. To avoid such redundancies, minimal nested parametrizations are introduced, and we show to what extent integrating this nested structure in the Gaussian process formalism can increase the robustness of the results, make easier the error control, and limit as much as possible over-fitting.

The outline of this chapter is as follows. First, Section 3.2 presents the theoretical framework for the definition of a Gaussian-process regression with a linear polynomial trend. Then, the nested polynomial trends we propose are detailed in Section 3.3. At last, the efficiency of the method is illustrated on a series of analytic examples in Section 3.4.

3.2 Gaussian process predictors

3.2.1 General framework

For $d \geq 1$, let $L^2(\mathbb{X}, \mathbb{R})$ be the space of square integrable functions on any compact subset \mathbb{X} of \mathbb{R}^d , with values in \mathbb{R} , equipped with the inner product $(\cdot, \cdot)_{\mathbb{X}}$, and the associated norm $\|\cdot\|_{\mathbb{X}}$, such that for all u and v in $L^2(\mathbb{X}, \mathbb{R})$,

$$(u, v)_{\mathbb{X}} := \int_{\mathbb{X}} u(\mathbf{x})v(\mathbf{x})d\mathbf{x}, \quad \|u\|_{\mathbb{X}}^2 := (u, u)_{\mathbb{X}}. \quad (3.2.1)$$

If \mathbb{X} is not compact, it is possible to introduce a weighted L_2 space.

Let \mathcal{S} be a physical system, whose response depends on a d -dimensional input vector $\mathbf{x} = (x_1, \dots, x_d)$, and whose performance can be evaluated from the computation of a quantity of interest, $y(\mathbf{x})$. Function y is a deterministic mapping that is assumed to be an element of $L^2(\mathbb{X}, \mathbb{R})$. In this chapter, we suppose that the maximal available information about y is a set of n code evaluations at the points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ in \mathbb{X} . Given this information, we are interested in the identification of the *best* predictor \hat{y} of y .

In that context, the Gaussian process regression (GPR), or Kriging, plays a major role [Sacks et al., 1989; Oakley and O'Hagan, 2002; Santner et al., 2003; Rasmussen and Williams, 2006]. It is indeed able to provide a prediction of $y(\mathbf{x})$, which is optimal in the class of the linear predictors of y , and whose precision can be *a posteriori* quantified. Such a method considers function y as a sample path of a real-valued Gaussian stochastic process Y . Let μ and C be respectively the mean and the covariance functions of Y :

$$Y(\cdot) \sim \text{GP}(\mu(\cdot), C(\cdot, \cdot)). \quad (3.2.2)$$

Besides, a set of observations of y is available. These observations are gathered in a n -dimensional vector:

$$\mathbf{y}^{\text{obs}} = \left(y^{(1)} = y(\mathbf{x}^{(1)}), \dots, y^{(n)} = y(\mathbf{x}^{(n)}) \right), \quad (3.2.3)$$

such that $\mathbb{P}(\cdot | \mathbf{y}^{\text{obs}})$ and $\mathbb{E}[\cdot | \mathbf{y}^{\text{obs}}]$ denote the conditional probability and conditional mathematical expectation respectively.

Therefore, gathering in the vector $\boldsymbol{\mu}$ and in the matrix \mathbf{R} the evaluations of μ and C at the available points, such that:

$$\begin{cases} \boldsymbol{\mu} := \left(\mu(\mathbf{x}^{(1)}), \dots, \mu(\mathbf{x}^{(n)}) \right), \\ \mathbf{R}_{ij} := C(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}), \quad 1 \leq i, j \leq n, \end{cases} \quad (3.2.4)$$

it can be shown [O'Hagan, 1978] that if matrix \mathbf{R} is invertible, then:

$$Y(\cdot) | \mu, C, \mathbf{y}^{\text{obs}} \sim \text{GP}(\mu^c(\cdot), C^c(\cdot, \cdot)), \quad (3.2.5)$$

where, for all \mathbf{x}, \mathbf{x}' in \mathbb{X} :

$$\begin{cases} \mu^c(\mathbf{x}) := \mu(\mathbf{x}) + \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \\ C^c(\mathbf{x}, \mathbf{x}') := C(\mathbf{x}, \mathbf{x}') - \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1} \mathbf{r}(\mathbf{x}'), \\ \mathbf{r}(\mathbf{x}) := \left(C(\mathbf{x}, \mathbf{x}^{(1)}), \dots, C(\mathbf{x}, \mathbf{x}^{(n)}) \right). \end{cases} \quad (3.2.6)$$

Under this formalism, also known as simple Kriging (see Section 1.4), the best prediction of y in an unobserved point \mathbf{x} is given by the mean value of $(Y(\mathbf{x}) | \mathbf{y}^{\text{obs}})$, $\mu^c(\mathbf{x})$, whereas $C^c(\mathbf{x}, \mathbf{x})$ quantifies the trust we can put in that prediction.

In practice, it appears that \mathbf{R} may not be invertible due to numerical reasons. This can generally be overcome by adding a small nugget to the covariance matrix and optimizing with respect to it too (see [Gramacy and Lee, 2012]).

3.2.2 Choice of the covariance function

Without information about the regularity of y , function C is generally chosen in general parametric families. In this chapter, function C is supposed to be an element of the Matern-5/2 class, such that for all \mathbf{x}, \mathbf{x}' in \mathbb{X} :

$$C(\mathbf{x}, \mathbf{x}') := \sigma^2 \prod_{i=1}^d (1 + \sqrt{5}h_i + 5h_i^2/3) \exp(-\sqrt{5}h_i), \quad h_i = |x_i - x'_i|/\ell_i. \quad (3.2.7)$$

Hence, covariance function C is characterized by a vector of hyper-parameters, $\Theta := (\sigma, \ell_1, \dots, \ell_d)$, whose values also have to be conditioned by \mathbf{y}^{obs} . More details about other usual parametric expressions for C can be found in Santner et al. [2003]. A *full Bayesian* approach would then require the introduction of a prior distribution for this vector, and the use of sampling techniques (such as Monte Carlo Markov Chains [Rubinstein and Kroese, 2008]) to approximate the posterior distribution of $(Y \mid \mathbf{y}^{\text{obs}})$ [Handcock and Stein, 1993; Kennedy and O'Hagan, 2001; Bilonis et al., 2013]. In this chapter, we will adopt an alternative approach, which consists in conditioning all the results by the maximum likelihood estimate of the covariance parameters. This method, which is generally called *plug-in* approach, has been used in many papers for the definition of Gaussian process-based predictors, as it presents a good compromise between complexity, efficiency, and errors control [Bichon et al., 2008; Bect et al., 2012]. In that case, explicit formula can be derived to evaluate the relevance of the GPR-based metamodel from a cross validation procedure [Dubrule, 1983].

3.2.3 Choice of the mean function

In the same way as for the covariance function, the mean function of Y is supposed to be parametrized by a p -dimensional vector β . In the general case, the computation of $\mathbb{E}[Y(\mathbf{x}) \mid \mathbf{y}^{\text{obs}}]$ is not direct, but if:

- covariance function C is known,
- μ is linear with respect to β , that is to say it exists a p -dimensional vector-valued function \mathbf{h} such that $\mu(\mathbf{x}) = \mathbf{h}(\mathbf{x})^T \beta$,
- β is uniformly distributed on \mathbb{R}^p (improper prior distribution),

then a Universal Kriging predictor can be defined (see Section 1.4 for further details):

$$Y(\cdot) \mid C, \mathbf{y}^{\text{obs}} \sim \text{GP}(\mu^c(\cdot), C^c(\cdot, \cdot)), \quad (3.2.8)$$

$$\left\{ \begin{array}{l} \mu^c(\mathbf{x}) := \mathbf{h}(\mathbf{x})^T \widehat{\beta} + \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{H} \widehat{\beta}), \\ C^c(\mathbf{x}, \mathbf{x}') := C(\mathbf{x}, \mathbf{x}') - \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1} \mathbf{r}(\mathbf{x}') + \mathbf{u}(\mathbf{x})^T (\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{u}(\mathbf{x}'), \\ \widehat{\beta} := (\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}^{-1} \mathbf{y}^{\text{obs}}, \\ \mathbf{u}(\mathbf{x}) := \mathbf{H}^T \mathbf{R}^{-1} \mathbf{r}(\mathbf{x}) - \mathbf{h}(\mathbf{x}), \\ \mathbf{H}_{ij} := f_j(\mathbf{x}^{(i)}), \quad 1 \leq j \leq p, 1 \leq i \leq n, \end{array} \right. \quad (3.2.9)$$

where the term $\mathbf{u}(\mathbf{x})^T (\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{u}(\mathbf{x}')$ can be interpreted as the prediction uncertainty that is due to the estimation of $\boldsymbol{\beta}$. Under these assumptions, the best prediction of $y(\mathbf{x})$ is now given by μ^c . The last thing that can be done to minimize $\|y - \mu^c\|_{\mathbb{X}}$ is working on the choice of \mathbf{h} .

Without information about y , polynomials are generally chosen for \mathbf{h} . Indeed, the set $\{m_{\boldsymbol{\alpha}}, \boldsymbol{\alpha} \in \mathbb{N}^d\}$, with

$$m_{\boldsymbol{\alpha}}(\mathbf{x}) := x_1^{\alpha_1} \times \cdots \times x_d^{\alpha_d}, \quad \mathbf{x} \in \mathbb{X}, \quad (3.2.10)$$

defines a basis of $L^2(\mathbb{X}, \mathbb{R})$. For a given value of p , characterizing \mathbf{h} amounts at identifying the best p -dimensional subset of $\{m_{\boldsymbol{\alpha}}, \boldsymbol{\alpha} \in \mathbb{N}^d\}$ to minimize $\|y - \mu^c\|_{\mathbb{X}}$.

In practice, this optimization problem over a very vast space is replaced by an optimization over a finite dimensional subset of $\{m_{\boldsymbol{\alpha}}, \boldsymbol{\alpha} \in \mathbb{N}^d\}$. Different truncation schemes have been proposed to choose such a relevant subset, which are mostly based on the assumption that the most influential elements of $\{m_{\boldsymbol{\alpha}}, \boldsymbol{\alpha} \in \mathbb{N}^d\}$ correspond to the elements of lowest total polynomial order. Denoting by r the maximal polynomial order of the projection basis, we can introduce:

$$\mathcal{P}(r, d) := \{m_{\boldsymbol{\alpha}} \mid \boldsymbol{\alpha} \in \mathbb{N}^d, \sum_{i=1}^d |\alpha_i| \leq r\}. \quad (3.2.11)$$

By construction, it can be noticed that the cardinal $\mathcal{C}(r, d)$ of $\mathcal{P}(r, d)$ increases exponentially with respect to r and d :

$$\mathcal{C}(r, d) = (d + r)! / (d! \times r!). \quad (3.2.12)$$

For $p \leq \mathcal{C}(r, d)$, vector \mathbf{h} can finally be searched using a penalization technique, such as the Least Angle Regression (LAR) method [Hastie et al., 2002; Efron et al., 2004; Blatman and Sudret, 2011], which allows disregarding insignificant terms. Such an approach will be referred as "LAR+UK" approach in the following.

3.3 Nested polynomial trends for Gaussian process predictors

As presented in Introduction, we are interested in identifying the best predictor of y in any unobserved point \mathbf{x} in \mathbb{X} , when the maximal information is a fixed number of code evaluations. Instead of considering sparse representations for the parametrization of the mean function in the GPR formalism, this section proposes to focus on nested polynomial representations. First, the notations and the motivations for this new parametrization are presented. Then, it is explained why and how it is integrated in the GPR formalism. Finally, a method to *a posteriori* evaluate the projection error is introduced.

3.3.1 Nested polynomial representations

Using the notations given by Eqs. (3.2.11) and (3.2.12), for p_2, p_1, d_2 in \mathbb{N}^* , let $\mathbf{m}^{(p_2, u_2)}$ and $\mathbf{m}^{(p_1, u_1)}$ be the vector-valued functions that gather all the elements of $\mathcal{P}(p_2, u_2)$ and $\mathcal{P}(p_1, u_1)$ respectively, and let $\mathcal{C}(p_2, u_2)$ and $\mathcal{C}(p_1, u_1)$ be their respective dimensions. The elements of these two vectors are sorted in an increasing total polynomial order. In particular, it comes:

$$m_1^{(p_2, u_2)} = m_1^{(p_1, u_1)} = 1, \quad (3.3.1)$$

where $u_1 = d$.

Hence, for all $(u_2 \times \mathcal{C}(p_1, u_1))$ -dimensional matrix \mathbf{A} and all $\mathcal{C}(p_2, u_2)$ -dimensional vector β_2 , the mapping

$$\mathbf{x} \mapsto \mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}) \quad (3.3.2)$$

is a function with values in \mathbb{R}^{u_2} , and the mapping

$$\mathbf{x} \mapsto \mathbf{m}^{(p_2, u_2)}(\mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}))^T \beta_2 \quad (3.3.3)$$

defines a nested polynomial representation. For $u_1 = d > 1$, such a representation allows us to model separately the dependence structure between the different input parameters, which is characterized by p_2 and u_2 , and the individual actions of each input parameter, which are characterized by the polynomial order p_1 (considering different values of p_1 for each input could eventually be done to optimize such a two-scale modeling). Hence, analyzing the optimal values of p_2 , u_2 and p_1 can bring information about the structure of y . For instance, if $p_2 = 1$ and $u_2 = d$, then y is just an additive model, up to a transformation of its input parameters. In the same manner, a value of p_1 strictly greater than 1 tends to say that the relation between \mathbf{x} and y is multi-scale.

Another interesting property of this nested structure comes from the fact that, for all \mathbf{x} in \mathbb{R}^d :

$$\begin{aligned} \mathbf{m}^{(p_2, u_2)}(\mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}))^T \beta_2 &= \sum_{0 \leq |\alpha_1| + \dots + |\alpha_{u_2}| \leq p_2} (\beta_2)_{(\alpha_1, \dots, \alpha_{u_2})} \times \prod_{i=1}^{u_2} \left(\sum_{k=1}^{\mathcal{C}(p_1, u_1)} \mathbf{A}_{ik} m_k^{(p_1, u_1)}(\mathbf{x}) \right)^{\alpha_i}, \\ &= \sum_{0 \leq |\tilde{\alpha}_1| + \dots + |\tilde{\alpha}_{u_1}| \leq p_2 \times p_1} x_1^{\tilde{\alpha}_1} \times \dots \times x_{u_1}^{\tilde{\alpha}_{u_1}} \tilde{c}_{\tilde{\alpha}}(\mathbf{A}, \beta_2; u_2), \end{aligned} \quad (3.3.4)$$

where $\tilde{c}_{\tilde{\alpha}}(\mathbf{A}, \beta_2; u_2)$ is the projection coefficient of $\mathbf{m}^{(p_2, u_2)}(\mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}))^T \beta_2$ on $x_1^{\tilde{\alpha}_1} \times \dots \times x_{u_1}^{\tilde{\alpha}_{u_1}}$. Hence, function $\mathbf{x} \mapsto \mathbf{m}^{(p_2, u_2)}(\mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}))^T \beta_2$ is in $\text{Span}\{\mathcal{P}(p_2 \times p_1, u_1)\}$, while being characterized by only $\mathcal{C}(p_2, u_2) + u_2 \times \mathcal{C}(p_1, u_1)$ parameters. Thus, by choosing u_2 such that the ratio $(\mathcal{C}(p_2, u_2) + u_2 \times \mathcal{C}(p_1, u_1)) / \mathcal{C}(p_2 \times p_1, u_1)$ is small, it is possible to parametrize polynomial families with very high cardinality, with only a reduced number of parameters. Such a parametrization is however redundant, in the sense that several distinct values of \mathbf{A} and β_2 lead to the same nested representations. From Eq. (3.3.4), it can be seen that some of these redundancies can be avoided by imposing that:

$$\begin{cases} \mathbf{A}_{i1} = 0, \\ \sum_{k=1}^{\mathcal{C}(p_1, u_1)} \mathbf{A}_{ik}^2 = 1, \end{cases} \quad 1 \leq i \leq u_2. \quad (3.3.5)$$

For fixed values of p_2 and p_1 , it is clear that ratio $(\mathcal{C}(p_2, u_2) + u_2 \times \mathcal{C}(p_1, u_1)) / \mathcal{C}(p_2 \times p_1, u_1)$ is minimal when $u_2 = 1$. However, considering higher values of u_2 strongly increases the flexibility of the nested representation to approximate function y . In this chapter, as a compromise between flexibility and minimal parametrization, for all $2 \leq k \leq \mathcal{C}(p_1, u_1)$, we thus propose to fix to zero all the components of $(\mathbf{A}_{1k}, \dots, \mathbf{A}_{u_2k})$ but one. This means that each component of vector $\mathbf{m}^{(p_1, u_1)}(\mathbf{x})$ is used only once in the construction of $\mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x})$, and that only $\#\text{Coeff}(p_1, p_2, u_1, u_2) = \mathcal{C}(p_2, u_2) + (\mathcal{C}(p_1, u_1) - 1) - u_2$ independent parameters have to be

Values of d	$\mathcal{C}(p_2 \times p_1, u_1)$	$\#\text{Coeff}(p_1, p_2, u_1, u_2 = 1)$	$\#\text{Coeff}(p_1, p_2, u_1, u_2 = d)$
1	10	6	6
2	55	12	17
5	2002	58	106
10	92378	288	561
20	10015005	1773	3521

Table 3.1: Comparison between the dimension of the projection set, $\mathcal{C}(p_2 \times p_1, u_1)$, and the number of independent parameters to characterize the associated projection coefficients in the proposed nested approach, $\#\text{Coeff}(p_1, p_2, u_1, u_2) = \mathcal{C}(p_2, u_2) + (\mathcal{C}(p_1, u_1) - 1) - u_2$, for $p_1 = p_2 = 3$, $u_1 \in \{1, 2, 5, 10, 20\}$ and $u_2 \in \{1, d\}$.

fixed to span a $\mathcal{C}(p_2 \times p_1, u_1)$ -dimensional projection set. As it can be seen in Table 3.1 and as it will be shown in Section 3.4, this assumption is indeed very attractive in terms of dimension reduction while being particularly interesting for the modeling of complex phenomena with very limited information.

To simplify the notations of the next sections, these $\mathcal{C}(p_1, u_1) - 1$ non-zero coefficients of \mathbf{A} are supposed to be gathered in a vector $\boldsymbol{\beta}_1$, and we introduce the matrices $\mathbf{P}^{(p_1, u_1)}(\mathbf{x})$ such that for all $\mathbf{x} \in \mathbb{X}$:

$$\mathbf{P}^{(p_1, u_1)}(\mathbf{x})\boldsymbol{\beta}_1 := \mathbf{A}\mathbf{m}^{(p_1, u_1)}(\mathbf{x}). \quad (3.3.6)$$

For given values of $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$, we then denote by $\mu(\cdot; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$ the following nested representation:

$$\mu(\mathbf{x}; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2) := \mathbf{m}^{(p_2, u_2)}(\mathbf{P}^{(p_1, u_1)}(\mathbf{x})\boldsymbol{\beta}_1)^T \boldsymbol{\beta}_2, \quad \mathbf{x} \in \mathbb{X}. \quad (3.3.7)$$

Finally, for given values of u_2 , p_2 , p_1 , the most appropriate nested representation to approximate function y is given by $\mu(\cdot; \boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2^*)$, where $(\boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2^*)$ is the solution of the following optimization problem:

$$(\boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2^*) := \arg \min_{(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \in \mathcal{S}^*} \|y - \mu(\cdot; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2)\|_{\mathbb{X}}^2, \quad (3.3.8)$$

and the admissible searching set, \mathcal{S}^* , is a subset of $\mathbb{R}^{\mathcal{C}(p_1, u_1) - 1} \times \mathbb{R}^{\mathcal{C}(p_2, u_2)}$ that takes into account the constraints on $\boldsymbol{\beta}_1$ defined by Eqs. (3.3.5) and (3.3.6).

Three main difficulties arise from the optimization problem defined by Eq. (3.3.8). First, as the maximal information about y is a n -dimensional set of evaluations, for given values of $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$, the norm $\|y - \mu(\cdot; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2)\|_{\mathbb{X}}^2$ has to be approximated. If the evaluation points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ are (more or less) uniformly distributed on \mathbb{X} , a (rather) good estimation of this norm is given by its least squares approximation,

$$\frac{1}{n} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \mu(\mathbf{x}^{(i)}; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \right)^2 = \frac{1}{n} \|\mathbf{y}^{\text{obs}} - \mathbf{M}(\boldsymbol{\beta}_1)\boldsymbol{\beta}_2\|^2, \quad (3.3.9)$$

where the vector \mathbf{y}^{obs} is defined by Eq. (3.2.3), and $\mathbf{M}(\boldsymbol{\beta}_1)$ is a $(n \times \mathcal{C}(p_2, u_2))$ -dimensional matrix such that:

$$(\mathbf{M}(\boldsymbol{\beta}_1))_{nk} = m_k^{(p_2, u_2)}(\mathbf{P}^{(p_1, u_1)}(\mathbf{x}^{(n)})\boldsymbol{\beta}_1), \quad 1 \leq n \leq n, \quad 1 \leq k \leq \mathcal{C}(p_2, u_2). \quad (3.3.10)$$

Noticing that for all (β_1, β_2) in \mathcal{S}^* ,

$$\left\| \mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1) (\mathbf{M}(\beta_1)^T \mathbf{M}(\beta_1))^{-1} \mathbf{M}(\beta_1)^T \mathbf{y}^{\text{obs}} \right\|^2 \leq \left\| \mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1) \beta_2 \right\|^2, \quad (3.3.11)$$

the solutions, β_1^* and β_2^* , of the minimization problem defined by Eq. (3.3.8) can respectively be approximated by the vectors β_1^{LS} and $\beta_2^{\text{LS}}(\beta_1^{\text{LS}})$, with:

$$\begin{cases} \beta_1^{\text{LS}} = \arg \min_{\beta_1 \in \mathcal{S}_{\beta_1}^*} \left\| \mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1) \beta_2^{\text{LS}}(\beta_1) \right\|^2, \\ \beta_2^{\text{LS}}(\beta_1) = (\mathbf{M}(\beta_1)^T \mathbf{M}(\beta_1))^{-1} \mathbf{M}(\beta_1)^T \mathbf{y}^{\text{obs}}, \end{cases} \quad (3.3.12)$$

where $\mathcal{S}_{\beta_1}^*$ is a subset of $\mathbb{R}^{C(p_1, u_1)-1}$ that also takes into account the constraints on β_1 defined by Eqs. (3.3.5) and (3.3.6).

The second difficulty comes from the fact that the minimization of the function $\beta_1 \mapsto \left\| \mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1) \beta_2^{\text{LS}}(\beta_1) \right\|^2$ can be complex. This is due to the fact that this mapping is strongly non-linear, leading to a strongly non-convex problem. For high values of p_2 , p_1 and u_2 , even if non-convex optimization algorithms such as simulated annealing or simplex algorithms [Brent, 1973] are used, there is no guarantee that the global minimum can be found in a reasonable computational time.

At last, there is a risk that $\left\| \mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1^{\text{LS}}) \beta_2^{\text{LS}}(\beta_1^{\text{LS}}) \right\|^2 / n$ strongly underestimates $\left\| \mathbf{y} - \mu(\cdot; \beta_1^{\text{LS}}, \beta_2^{\text{LS}}(\beta_1^{\text{LS}})) \right\|_{\mathbb{X}}^2$, as the same information is used twice: once for the optimization and once for the error estimation. To avoid such an over-fitting, classical Leave-One-Out (LOO) techniques (see Miller [1974]; Blatman and Sudret [2011]; Perrin et al. [2014]) have to be introduced to get a relevant approximation of $\left\| \mathbf{y} - \mu(\cdot; \beta_1^{\text{LS}}, \beta_2^{\text{LS}}(\beta_1^{\text{LS}})) \right\|_{\mathbb{X}}^2$.

3.3.2 Coupling nested representations and Gaussian processes

Once vector β_1^{LS} has been identified from the solving of Eq. (3.3.12), the notion of confidence intervals for the prediction of $y(\mathbf{x})$ at an unobserved point \mathbf{x} can be found back by assuming that y is a particular realization of a Gaussian stochastic process, whose statistical properties are given by:

$$Y(\cdot) \sim \text{GP} \left(\mu(\cdot; \beta_1^{\text{LS}}, \beta_2^{\text{LS}}(\beta_1^{\text{LS}})), C(\cdot, \cdot; \widehat{\Theta}^{\text{LS}}) \right), \quad (3.3.13)$$

where $\widehat{\Theta}^{\text{LS}}$ gathers the $d+1$ parameters of the Matern-5/2 covariance C defined by Eq. (3.2.7), which are solution of the following log-likelihood maximization problem:

$$\widehat{\Theta}^{\text{LS}} = \underset{\Theta \in (0, +\infty)^{d+1}}{\text{argmax}} - \frac{1}{2} \left[\begin{array}{c} n \log(2\pi) + \log(\det(\mathbf{R}(\Theta))) + \\ (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1^{\text{LS}}) \beta_2^{\text{LS}}(\beta_1^{\text{LS}}))^T \mathbf{R}(\Theta)^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1^{\text{LS}}) \beta_2^{\text{LS}}(\beta_1^{\text{LS}})) \end{array} \right]. \quad (3.3.14)$$

Such a naive coupling is nevertheless sub-optimal, as the values of β_1 and Θ are optimized separately: the nested structure does not take advantage of the Bayesian formalism, and reciprocally. Instead of such a two-steps approach, we propose in this chapter to directly adopt a Bayesian formalism for the estimation of β_1 and Θ . In the plug-in formalism, this means that the statistical properties of Y are now given by:

$$Y(\cdot) \sim \text{GP} \left(\mu(\cdot; \widehat{\beta}_1, \widehat{\beta}_2), C(\cdot, \cdot; \widehat{\Theta}) \right), \quad (3.3.15)$$

where $(\widehat{\beta}_1, \widehat{\beta}_2, \widehat{\Theta})$ is the solution of the following log-likelihood maximization problem:

$$(\widehat{\beta}_1, \widehat{\beta}_2, \widehat{\Theta}) = \underset{(\beta_1, \beta_2, \Theta) \in \mathcal{S}^{\text{adm}}}{\text{argmax}} - \frac{1}{2} \left[\begin{array}{c} n \log(2\pi) + \log(\det(\mathbf{R}(\Theta))) \\ + (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1)\beta_2)^T \mathbf{R}(\Theta)^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1)\beta_2) \end{array} \right], \quad (3.3.16)$$

where the admissible searching set, \mathcal{S}^{adm} , is a subset of $\mathbb{R}^{\mathcal{C}(p_1, u_1)-1} \times \mathbb{R}^{\mathcal{C}(p_2, u_2)} \times \mathbb{R}^{d+1}$ but is not trivial, as it first takes into account the constraints on β_1 defined by Eqs. (3.3.5) and (3.3.6), but also guarantees that $\mathbf{R}(\Theta)$ and $\mathbf{M}(\beta_1)^T \mathbf{R}(\Theta)^{-1} \mathbf{M}(\beta_1)$ are invertible.

For all $(\beta_1, \beta_2, \Theta)$ belonging to the admissible set, \mathcal{S}^{adm} , we denote by L the function such that:

$$L(\beta_1, \beta_2, \Theta) = \log(\det(\mathbf{R}(\Theta))) + (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1)\beta_2)^T \mathbf{R}(\Theta)^{-1} (\mathbf{y}^{\text{obs}} - \mathbf{M}(\beta_1)\beta_2). \quad (3.3.17)$$

It is interesting to notice that, in the same manner as in Section 3.3.1,

$$L(\beta_1, \beta_2^{\text{LS}}(\beta_1, \Theta), \Theta) \leq L(\beta_1, \beta_2, \Theta), \quad (3.3.18)$$

$$\beta_2^{\text{LS}}(\beta_1, \Theta) := (\mathbf{M}(\beta_1)^T \mathbf{R}(\Theta)^{-1} \mathbf{M}(\beta_1))^{-1} \mathbf{M}(\beta_1)^T \mathbf{R}(\Theta) \mathbf{y}^{\text{obs}}. \quad (3.3.19)$$

It comes:

$$\begin{cases} (\widehat{\beta}_1, \widehat{\Theta}) = \arg \min_{(\beta_1, \Theta)} \mathcal{L}(\beta_1, \Theta), \\ \widehat{\beta}_2 = \left(\mathbf{M}(\widehat{\beta}_1)^T \mathbf{R}(\widehat{\Theta})^{-1} \mathbf{M}(\widehat{\beta}_1) \right)^{-1} \mathbf{M}(\widehat{\beta}_1)^T \mathbf{R}(\widehat{\Theta}) \mathbf{y}^{\text{obs}}, \end{cases} \quad (3.3.20)$$

where:

$$\mathcal{L}(\beta_1, \Theta) := L(\beta_1, \beta_2^{\text{LS}}(\beta_1, \Theta), \Theta). \quad (3.3.21)$$

Function $(\beta_1, \Theta) \mapsto \mathcal{L}(\beta_1, \Theta)$ being strongly non-regular and non-convex, it is proposed to work iteratively on the values of β_1 and Θ . Two reasons motivate this separation. First, the actions of β_1 and Θ on $\mathcal{L}(\beta_1, \Theta)$ being very different, dividing the optimization problem tends to regularize the mappings on which the minimization is carried out. Second, by reducing each searching set, each minimization is made easier. Therefore, for a given convergence tolerance ε , Algorithm 1 is introduced for the minimization of \mathcal{L} . The convergence of such an iterative algorithm to the global minimum of \mathcal{L} is of course not guaranteed, but it appeared on a series of numerical examples that it allowed us to identify good approximations of $(\widehat{\beta}_1, \widehat{\Theta})$ at a reasonable computational cost. As the minimization problem defined by Eq. (3.3.20) is not convex, better approximations of $\widehat{\beta}_1$ can be obtained by repeating several times Algorithm 1, with random initialization of vectors $(\beta_1)_0$ in $\mathcal{S}_{\beta_1}^*$.

3.3.3 Linearization of the nested polynomial trend

Even for small values of p_2 , p_1 and u_2 , the quantity $\mathcal{L}(\beta_1, \Theta)$ is sensitive to small changes in the values of β_1 and Θ , which makes the solving of the optimization problem defined by Eq. (3.3.20) difficult. In that context, it can be interesting to linearize the nested polynomial trend around the solutions given by Algorithm 1, $\widehat{\beta}_1$ and $\widehat{\beta}_2$, and then work on the compensations $(\beta_1 - \widehat{\beta}_1)$ and $(\beta_2 - \widehat{\beta}_2)$ that could make the prediction of function y better. In the vicinity of $\widehat{\beta}_1$ and $\widehat{\beta}_2$, for all \mathbf{x} in \mathbb{X} , it comes:

$$\mu(\mathbf{x}; \beta_1, \beta_2) \approx \left(\mathbf{h}_1(\mathbf{x}; \widehat{\beta}_1, \widehat{\beta}_2), \mathbf{h}_2(\mathbf{x}; \widehat{\beta}_1) \right)^T (\beta_1 - \widehat{\beta}_1, \beta_2), \quad (3.3.22)$$

1 Initialization: $L_1 = 0, L_2 = +\infty, \beta_1^* = (\beta_1)_0 \in \mathcal{S}_{\beta_1}^*$;
2 while $|L_2 - L_1| > \varepsilon$ **do**
3 $L_1 = L_2$;
4 $\Theta^* = \arg \max_{\Theta} \mathcal{L}(\beta_1^*, \Theta)$;
5 $\beta_1^* = \arg \max_{\beta_1} \mathcal{L}(\beta_1, \Theta^*)$;
6 $L_2 = \min(L_2, \mathcal{L}(\beta_1^*, \Theta^*))$;
7 end
8 $\hat{\beta}_1 \approx \beta_1^*, \hat{\Theta} \approx \Theta^*$.

Algorithm 1: Iterative minimization of function \mathcal{L} .

$$\mathbf{h}_1(\mathbf{x}; \hat{\beta}_1, \hat{\beta}_2) = \mathbf{P}^{(p_1, u_1)}(\mathbf{x})^T \mathbf{D}(\mathbf{P}^{(p_1, u_1)}(\mathbf{x}) \hat{\beta}_1)^T \hat{\beta}_2, \quad (3.3.23)$$

$$\mathbf{h}_2(\mathbf{x}; \hat{\beta}_1) = \mathbf{m}^{(p_2, u_2)}(\mathbf{P}^{(p_1, u_1)}(\mathbf{x}) \hat{\beta}_1), \quad (3.3.24)$$

$$(\mathbf{D}(\mathbf{z}))_{kj} := \frac{\partial m_k^{(p_2, u_2)}}{\partial z_j}(\mathbf{z}), \quad 1 \leq j \leq u_2, \quad 1 \leq k \leq \mathcal{C}(p_2, u_2), \quad \mathbf{z} \in \mathbb{R}^{u_2}. \quad (3.3.25)$$

Now, let us denote by $\beta := (\beta_1 - \hat{\beta}_1, \beta_2)$ the new vector of parameters we need to determine, and by $\mathbf{h} := (\mathbf{h}_1(\cdot; \hat{\beta}_1, \hat{\beta}_2), \mathbf{h}_2(\cdot; \hat{\beta}_1))$ the new set of projection functions. Conditioned by the values of $\hat{\beta}_1, \hat{\beta}_2$ and $\hat{\Theta}$, the formalism introduced in Section 3.2.3 is found back:

$$Y(\cdot) \sim \text{GP} \left(\mathbf{h}(\cdot)^T \beta, C(\cdot, \cdot) \right), \quad (3.3.26)$$

such that the distribution of $(Y | \mathbf{y}^{\text{obs}})$ can be calculated analytically. Its mean value can directly be used to predict the values of y , and its covariance function can allow us to quantify the confidence we can put in these predictions.

We underline at least two advantages for the linearization. First, the distribution of $(Y | \mathbf{y}^{\text{obs}})$ will be less dependent on the convergence properties of Algorithm 1, which are not easy to control. Secondly, as the covariance function of $(Y | \mathbf{y}^{\text{obs}})$ integrates the uncertainty associated with the least squares estimation of β , that is to say the uncertainty associated with the estimation of β_1 and β_2 in the vicinity of $\hat{\beta}_1$ and $\hat{\beta}_2$, the confidence intervals associated with these predictions are expected to be more adapted.

3.3.4 Error evaluation

According to the previous Sections and to Eq. (3.2.9), for given values of truncation parameters p_2, p_1 and u_2 , we propose to use the deterministic function $\hat{y}^{\text{nest}}(\mathbf{x})$, such that:

$$\hat{y}^{\text{nest}}(\mathbf{x}) = \mathbf{h}(\mathbf{x}; \hat{\beta}_1, \hat{\Theta})^T \hat{\beta}(\hat{\beta}_1, \hat{\Theta}) + \mathbf{r}(\mathbf{x}; \hat{\Theta})^T \mathbf{R}(\hat{\Theta})^{-1} \left(\mathbf{y}^{\text{obs}} - \mathbf{H}(\hat{\beta}_1, \hat{\Theta}) \hat{\beta}(\hat{\beta}_1, \hat{\Theta}) \right), \quad (3.3.27)$$

$$\hat{\beta}(\hat{\beta}_1, \hat{\Theta}) := (\mathbf{H}(\hat{\beta}_1, \hat{\Theta})^T \mathbf{R}(\hat{\Theta})^{-1} \mathbf{H}(\hat{\beta}_1, \hat{\Theta}))^{-1} \mathbf{H}(\hat{\beta}_1, \hat{\Theta})^T \mathbf{R}(\hat{\Theta})^{-1} \mathbf{y}^{\text{obs}}, \quad (3.3.28)$$

to predict the value of $y(\mathbf{x})$ for all \mathbf{x} in \mathbb{X} , where:

- vectors $\hat{\beta}_1$ and $\hat{\Theta}$ are the solutions of the optimization problem given by Eq. (3.3.20), under the additional condition that the matrix $\mathbf{H}(\hat{\beta}_1, \hat{\Theta})^T \mathbf{R}(\hat{\Theta})^{-1} \mathbf{H}(\hat{\beta}_1, \hat{\Theta})$ is invertible,

- vector \mathbf{y}^{obs} is defined by Eq. (3.2.3),
- the function $\mathbf{x} \mapsto \mathbf{h}(\mathbf{x}; \hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})$ gathers the most influential terms of the vector-valued function $\left(\mathbf{h}_1(\cdot; \hat{\boldsymbol{\beta}}_1, \boldsymbol{\beta}_2^{\text{LS}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})), \mathbf{h}_2(\cdot; \hat{\boldsymbol{\beta}}_1)\right)$, which have been identified from a LAR procedure,
- $\mathbf{H}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) := [\mathbf{h}(\mathbf{x}^{(1)}; \hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) \cdots \mathbf{h}(\mathbf{x}^{(n)}; \hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})]$ is the matrix that gathers the evaluations of $\mathbf{h}(\cdot; \hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})$ at the available code evaluations,
- and for all $1 \leq i, j \leq n$, $\mathbf{R}(\hat{\boldsymbol{\Theta}})_{ij} = C(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ and $r_i(\mathbf{x}; \hat{\boldsymbol{\Theta}}) = C(\mathbf{x}, \mathbf{x}^{(i)})$, with C the Matern-5/2 covariance function of parameters $\hat{\boldsymbol{\Theta}}$.

In the same manner as in Section 3.2, when function y is only known through a limited number of evaluations, classical Leave-One-Out (LOO) techniques have to be introduced to approximate the relevance of such a predictor:

$$\|y - \hat{y}^{\text{nest}}\|_{L_2}^2 \approx \epsilon_{\text{LOO}}^2 := \frac{1}{n} \sum_{i=1}^n \left(y(\mathbf{x}^{(i)}) - \hat{y}_{-i}^{\text{nest}}(\mathbf{x}^{(i)}) \right)^2, \quad (3.3.29)$$

where, for all $1 \leq i \leq n$, the function $\hat{y}_{-i}^{\text{nest}}$ has been constructed in the same manner as \hat{y}^{nest} , but using the $n - 1$ evaluations of the code in $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i-1)}, \mathbf{x}^{(i+1)}, \dots, \mathbf{x}^{(n)}\}$ only.

In order to reduce the computational cost associated with the evaluation of ϵ_{LOO}^2 , it is interesting to notice (see Dubrule [1983] for further details) that, for all $1 \leq i \leq n$:

$$y(\mathbf{x}^{(i)}) - \hat{y}_{-i}^{\text{nest}}(\mathbf{x}^{(i)}) = \frac{(\hat{\mathbf{C}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) \mathbf{y}^{\text{obs}})_i}{\hat{\mathbf{C}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})_{ii}}, \quad (3.3.30)$$

$$\begin{aligned} \hat{\mathbf{C}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) &= \mathbf{R}(\hat{\boldsymbol{\Theta}})^{-1} - \\ &\mathbf{R}(\hat{\boldsymbol{\Theta}})^{-1} \mathbf{H}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) \left(\mathbf{H}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})^T \mathbf{R}(\hat{\boldsymbol{\Theta}})^{-1} \mathbf{H}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) \right)^{-1} \mathbf{H}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})^T \mathbf{R}(\hat{\boldsymbol{\Theta}})^{-1}. \end{aligned} \quad (3.3.31)$$

LOO error ϵ_{LOO}^2 can then be approximated by:

$$\epsilon_{\text{LOO}}^2 \approx \hat{\epsilon}_{\text{LOO}}^2 := \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2, \quad \hat{e}_i^2 := \left[\frac{(\hat{\mathbf{C}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}}) \mathbf{y}^{\text{obs}})_i}{\hat{\mathbf{C}}(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\Theta}})_{ii}} \right]^2. \quad (3.3.32)$$

Such an approximation is however conditioned by the values of $\hat{\boldsymbol{\beta}}_1$ and $\hat{\boldsymbol{\Theta}}$, which are computed using all the code evaluations. In order to be more precise, it can be noticed that for all $\boldsymbol{\beta}_1, \boldsymbol{\Theta}$, $1 \leq i \leq n$:

$$\mathcal{L}(\boldsymbol{\beta}_1, \boldsymbol{\Theta}) = \mathcal{L}_{-i}(\boldsymbol{\beta}_1, \boldsymbol{\Theta}) + \frac{(\tilde{\mathbf{C}}(\boldsymbol{\beta}_1, \boldsymbol{\Theta}) \mathbf{y}^{\text{obs}})_i^2}{\tilde{\mathbf{C}}(\boldsymbol{\beta}_1, \boldsymbol{\Theta})_{ii}}, \quad (3.3.33)$$

$$\tilde{\mathbf{C}}(\boldsymbol{\beta}_1, \boldsymbol{\Theta}) = \mathbf{R}(\boldsymbol{\Theta})^{-1} \{ \mathbf{I} - \mathbf{M}(\boldsymbol{\beta}_1) (\mathbf{M}(\boldsymbol{\beta}_1)^T \mathbf{R}(\boldsymbol{\Theta})^{-1} \mathbf{M}(\boldsymbol{\beta}_1))^{-1} \mathbf{M}(\boldsymbol{\beta}_1)^T \mathbf{R}(\boldsymbol{\Theta})^{-1} \}, \quad (3.3.34)$$

where \mathbf{I} is the identity matrix and $\mathcal{L}_{-i}(\boldsymbol{\beta}_1, \boldsymbol{\Theta})$ is the evaluation of function $\mathcal{L}(\boldsymbol{\beta}_1, \boldsymbol{\Theta})$ based on the $n - 1$ evaluations of the code in $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i-1)}, \mathbf{x}^{(i+1)}, \dots, \mathbf{x}^{(n)}\}$ only. Hence, in the optimization process leading us to the identification of $\hat{\boldsymbol{\beta}}_1$ and $\hat{\boldsymbol{\Theta}}$, let $\{((\boldsymbol{\beta}_1)_i, \boldsymbol{\Theta}_i), 1 \leq i \leq n_{\text{test}}\}$ be the n_{test} values of $\boldsymbol{\beta}_1$ and $\boldsymbol{\Theta}$, in which function \mathcal{L} has been evaluated. With a

very limited additional computational cost, we can then define, for all $1 \leq i \leq n$, the LOO evaluations of $\widehat{\beta}_1$ and $\widehat{\Theta}$, which are denoted by $\left(\widehat{\beta}_1\right)_{-i}$ and $\widehat{\Theta}_{-i}$ respectively, and which are given by:

$$\left(\widehat{\beta}_1\right)_{-i}, \widehat{\Theta}_{-i} = \arg \min_{(\beta_1, \Theta) \in \{((\beta_1)_i, \Theta_i), 1 \leq i \leq n_{\text{test}}\}} \mathcal{L}_{-i}(\beta_1, \Theta). \quad (3.3.35)$$

Finally, we can introduce error $\tilde{\epsilon}_{\text{LOO}}$, such that:

$$\|y - \widehat{y}^{\text{nest}}\|_{L_2}^2 \approx \tilde{\epsilon}_{\text{LOO}}^2 := \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2, \quad \tilde{e}_i^2 := \left[\frac{\left(\widehat{C}\left(\left(\widehat{\beta}_1\right)_{-i}, \widehat{\Theta}_{-i}\right) \mathbf{y}^{\text{obs}}\right)_i}{\widehat{C}\left(\left(\widehat{\beta}_1\right)_{-i}, \widehat{\Theta}_{-i}\right)_{ii}} \right]^2. \quad (3.3.36)$$

3.3.5 Convergence analysis

All the developments presented in Sections 3.3.1 and 3.3.2 are conditioned by the values of three truncation parameters, p_2 , p_1 and u_2 , which have to be identified from a convergence analysis. As presented in Section 3.3.1, we remind that the roles of p_2 , p_1 and u_2 in the modeling of y are different. Whereas p_2 and u_2 are associated with the modeling of the dependency structure between the input parameters, p_1 is associated with the individual transformation of each input. As a consequence, p_1 is strongly dependent on the dimension of vector β_1 , which parametrizes these individual transformations. On the contrary, this dimension of β_1 , which is equal to $\mathcal{C}(p_1, u_1) - 1 - u_2$, does not depend on p_2 , but depends only linearly on u_2 . Hence, increasing the values of p_2 and u_2 does not really increase the dimension of the search set for the identification of $\widehat{\beta}_1$, but makes the relation between β_1 and $\mathcal{L}(\beta_1, \Theta)$ much more complex.

For the choice of u_2 , p_2 and p_1 , maximal values u_2^{max} , p_2^{max} and q^{max} are *a priori* chosen. In this chapter, since we want to reduce the number of parameters on which the polynomial trend is based, only values of u_2 that are lower than d are considered: $u_2^{\text{max}} = d$. Finally, the optimal value of (u_2, p_1, p_2) is the one that gives the minimum LOO error among all these tested combinations of values:

$$(u_2^*, p_1^*, p_2^*) := \underset{\substack{1 \leq u_2 \leq d, \\ 1 \leq p_2 \leq p_2^{\text{max}}, \\ 1 \leq p_1 \leq p_1^{\text{max}}}}{\text{argmin}} \tilde{\epsilon}_{\text{LOO}}^2(u_2, p_1, p_2), \quad (3.3.37)$$

where error $\tilde{\epsilon}_{\text{LOO}}^2$ is defined by Eq. (3.3.32).

3.4 Applications

To illustrate the advantages of the nested structure presented in Section 3.3 for the modeling of the quantity of interest y , this section introduces a series of analytic examples, which are sorted with respect to the input set dimension, d . In each case, the proposed approach is compared to the "LAR+UK" approach, which has been described in Section 3.2. For each function y , let $\widehat{y}^{\text{nest}}$ and $\widehat{y}^{\text{LAR+UK}}$ be the best approximations of y we can get from the available information, when considering a nested polynomial trend and a simple polynomial trend, respectively. Let ϵ_{NEST}^2 and $\epsilon_{\text{LAR+UK}}^2$ be the associated normalized errors, such that:

$$\epsilon_{\text{NEST}}^2 = \|y - \widehat{y}^{\text{nest}}\|_{\mathbb{X}}^2 / \|g\|_{\mathbb{X}}^2, \quad (3.4.1)$$

$$\varepsilon_{\text{LAR+UK}}^2 = \|y - \hat{y}^{\text{LAR+UK}}\|_{\mathbb{X}}^2 / \|g\|_{\mathbb{X}}^2. \quad (3.4.2)$$

When dealing with a simple polynomial trend, it is reminded that the only truncation parameter that needs to be identified is the maximal total polynomial order, which will be denoted in the following by $p^{\text{LAR+UK}}$ for the sake of clarity. On the contrary, three truncation parameters have to be identified for the nested polynomial trends: p_2 , u_2 and p_1 . As a consequence, the required computational time to identify \hat{y}^{nest} can be much higher than the one required to identify $\hat{y}^{\text{LAR+UK}}$.

3.4.1 $d = 1$

In this part, we suppose that $d = 1$, and we fix $\mathbb{X} = [-1, 1]$. Three analytic expressions for y are then proposed:

- case 1: $y(x) = P_2 \circ P_1(x)$,
- case 2: $y(x) = \sin((x + 1)^3)$,
- case 3: $y(x) = \sin(20x) \cos(2x)$,

where, for all x in $[-1, 1]$:

$$\begin{cases} P_1(x) = \sum_{i=1}^5 c_i^{(1)} x^{i-1}, & \mathbf{c}^{(1)} = \frac{(0, -0.03, 0.5, -0.4, -0.5)}{\sqrt{0.03^2 + 0.5^2 + 0.4^2 + 0.5^2}}, \\ P_2(x) = \sum_{i=1}^5 c_i^{(2)} x^{i-1}, & \mathbf{c}^{(2)} = (-0.1, 0.2, 0.7, -0.2, -0.2). \end{cases} \quad (3.4.3)$$

The two first examples are based on chained codes. The third example is introduced to show that this nested structure for the mean can also be interesting for non-chained codes when few code evaluations are available.

For each case, Figure 3.1 compares the evolution of the errors $\varepsilon_{\text{NEST}}^2$ and $\varepsilon_{\text{LAR+UK}}^2$ with respect to n , the number of available evaluations of y . For each value of n , convergence analyses have been performed for both methods. The maximal values for the truncation parameters associated were fixed such that:

$$0 \leq p^{\text{LAR+UK}} \leq 20, \quad 0 \leq p_1, p_2 \leq 10, \quad u_2 = 1. \quad (3.4.4)$$

For the three applications, these convergence analyses lead us to relatively high values for these truncation parameters ($p_1 \geq 4$, $p_2 \geq 4$). As underlined in Section 3.3.1, this can be explained by the ability of the proposed nested structure to parametrize polynomial families with very high cardinality with only few parameters. This is particularly efficient when n is small compared to the number of oscillations of y .

In addition, Figure 3.2 compares the two approaches in terms of prediction for given values of n . In these figures we notice that the proposed method is particularly adapted to the cases when y presents a nested structure or is oscillating. This is particularly true when n is small compared to the complexity of y .

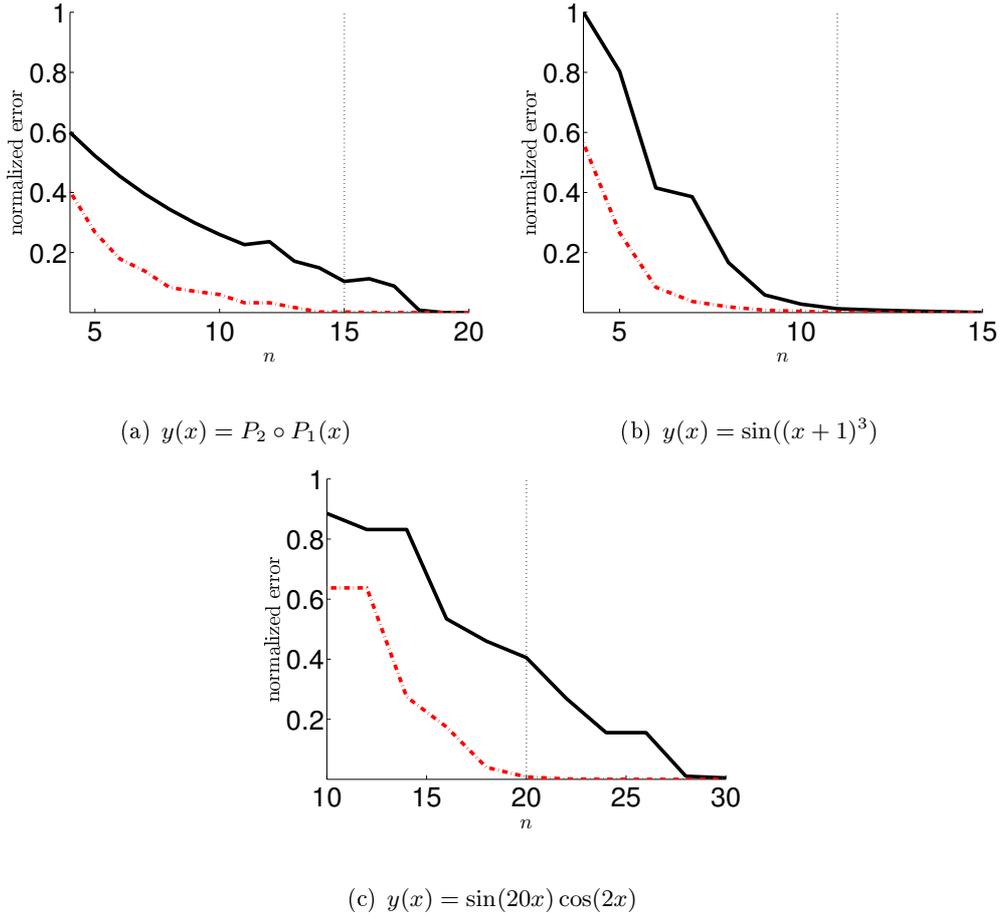


Figure 3.1: Evolution of the normalized L^2 errors with respect to n , the number of code evaluations. To be more representative, for each value of n , the LAR+UK and the proposed approaches have been repeated 10 times on randomly chosen learning sets. The curves correspond to the mean value of the errors associated with these 10 repetitions. Solid black line: evolution of the error associated with the LAR+UK approach, $\varepsilon_{\text{LAR+UK}}^2$. Red dotted line: evolution of the error associated with the proposed approach, $\varepsilon_{\text{NEST}}^2$. The vertical bar indicates moreover the value of n on which the results of Figure 3.2 are focused.

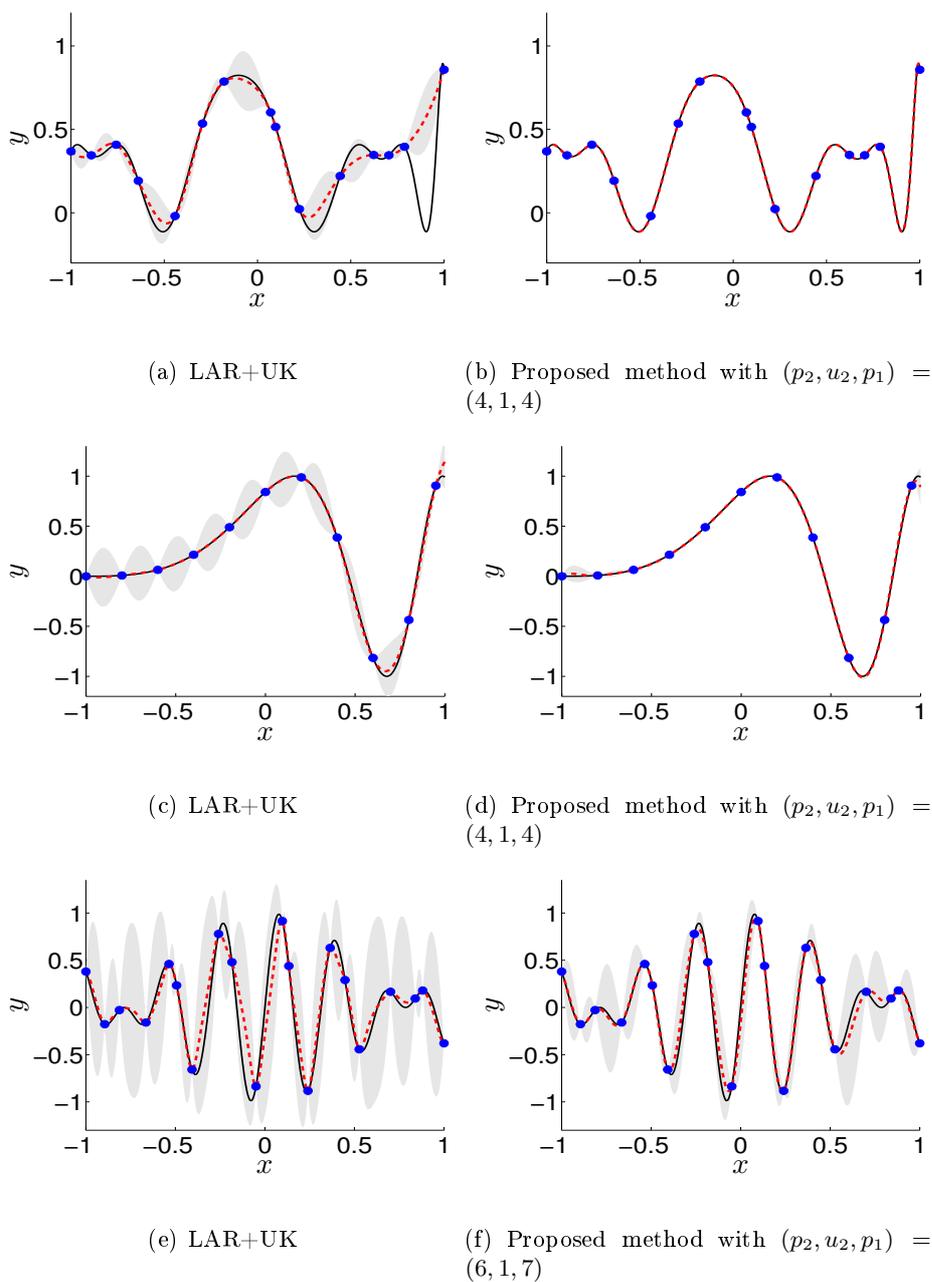


Figure 3.2: Efficiency of the proposed method to predict in an unobserved point the value of $y(x) = P_2 \circ P_1(x)$ with $n = 15$ (first row), $y(x) = \sin((x + 1)^3)$ with $n = 11$ (second row) and $y(x) = \sin(20x) \cos(2x)$ with $n = 20$ (third row). In each figure, the black solid line is the evolution of the quantity of interest, y , with respect to x , the blue points are the positions of the available observations of y , the red dotted line is the prediction of y based on an optimized LAR+UK approach (left column) or based on the proposed approach associated with optimized values of p_2 , u_2 and p_1 (right column). The grey areas correspond to the 95% confidence region for the prediction.

3.4.2 $d > 1$

The idea of this section is to show that the tendencies that were noticed in the one-dimensional cases are found back when considering multidimensional input spaces. To this end, let us consider the three following expressions of y , which can also be seen as particular chained codes, and the associated maximal values for the convergence analyses:

- Case 1: $d = 2$, $0 \leq p^{\text{LAR+UK}} \leq 20$, $0 \leq p_2 \leq 6$, $0 \leq p_1 \leq 10$, $1 \leq u_2 \leq d$.

$$g : \begin{cases} [-1, 1]^2 & \rightarrow \\ \mathbf{x} & \mapsto \end{cases} g^{2\text{D}}(\mathbf{x}) = (1 - x_1^2) \cos(7x_1) \times (1 - x_2^2) \sin(5x_2) \quad \cdot \quad (3.4.5)$$

- Case 2 (the Ishigami function): $d = 3$, $0 \leq p^{\text{LAR+UK}} \leq 20$, $0 \leq p_2 \leq 3$, $0 \leq p_1 \leq 10$, $1 \leq u_2 \leq d$.

$$g : \begin{cases} [-\pi, \pi]^3 & \rightarrow \\ \mathbf{x} = (x_1, x_2, x_3) & \mapsto \end{cases} g^{3\text{D}}(\mathbf{x}) = \sin(x_1) + 7 \sin(x_2)^2 + 0.1x_3^4 \sin(x_1) \quad \cdot \quad (3.4.6)$$

- Case 3: $d = 6$, $0 \leq p^{\text{LAR+UK}} \leq 10$, $0 \leq p_2 \leq 3$, $0 \leq p_1 \leq 10$, $1 \leq u_2 \leq d$.

$$g : \begin{cases} [-1, 1]^6 & \rightarrow \\ \mathbf{x} & \mapsto \end{cases} g^{6\text{D}}(\mathbf{x}) = g^{(1)} \circ \mathbf{g}^{(2)}(\mathbf{x}), \quad (3.4.7)$$

$$g^{(1)}(\mathbf{z}) = 0.1 \cos\left(\sum_{i=1}^6 z_i\right) + \sum_{i=1}^6 z_i^2, \quad \mathbf{z} \in \mathbb{R}^6, \quad (3.4.8)$$

$$\mathbf{g}^{(2)}(\mathbf{x}) = (\cos(\pi x_1 + 1), \cos(\pi x_2 + 2), \dots, \cos(\pi x_6 + 6)). \quad (3.4.9)$$

In the same manner as in Section 3.4.1, Figure 3.3 compares the evolution of errors $\varepsilon_{\text{NEST}}^2$ and $\varepsilon_{\text{LAR+UK}}^2$ with respect to n . As for the one-dimensional cases, it can be noticed in these figures that, for the studied examples, introducing a nested structure for the polynomial trend can allow us to make the L^2 error decrease by several orders of magnitude, especially when n is low. Moreover, these figures emphasize the interest of optimizing the values of the truncation parameter u_2 when dealing with multidimensional input spaces.

Note that for these examples, there is no information about the structure of the nested code. Adding some information about the relation between the inputs could be very useful to avoid testing too many values of p_1 , p_2 and u_2 .

As explained in Section 3.3.1, the values of p_2 , p_1 and u_2 that were obtained from the convergence analyses can give many information about the unknown structure of the quantity of interest. For the first example, the values $p_2 = 2$, $u_2 = 2$ and $p_1 > 2$ were most of the time chosen, which is coherent with the fact that $g^{2\text{D}}(x_1, x_2)$ is just the product of two functions that depend on x_1 and x_2 only. Hence, a particular attention has to be paid to the modeling of each input, rather than to the modeling of the dependence structure.

In the same manner, for the second example, most of the convergence analyses lead us to $u_2 = 3$ and $p_2 < p_1$, which also shows that the modeling of each input seems to be more important than the characterization of the relation between these modified inputs.

At last, for the third quantity of interest, which is a highly oscillating function in dimension $d = 6$, the convergence analyses seemed to encourage the values of p_2 and p_1 that lead to

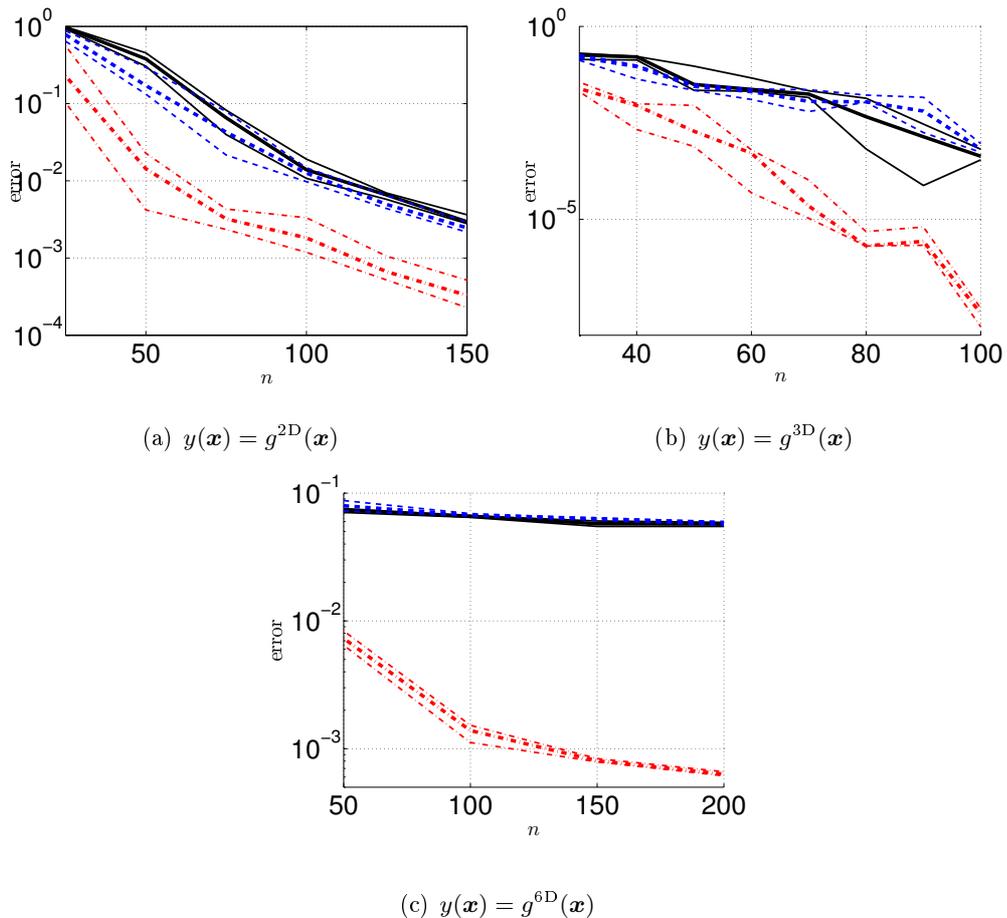


Figure 3.3: Evolution of the normalized L^2 errors with respect to n , the number of code evaluations. To be more representative, for each value of n , the LAR+UK and the proposed approaches have been repeated 10 times on randomly chosen learning sets. The curves correspond to the values of the 25% (thin line), the 50% (thick line) and the 75% (thin line) quantiles of the errors associated with these 10 repetitions. Solid black line: evolution of the error associated with the LAR+UK approach, $\varepsilon_{\text{LAR+UK}}^2$. Blue dotted line: evolution of the error associated with the proposed approach, $\varepsilon_{\text{NEST}}^2$, with $u_2 = 1$. Red dashed line: evolution of the error associated with the proposed approach, $\varepsilon_{\text{NEST}}^2$, with $1 \leq u_2 \leq d$.

the highest product $p_1 \times p_2$ (before over-fitting). This means that, for this example, it is interesting to approximate quantity of interest y by a complex polynomial representation that is characterized by a small number of parameters.

3.4.3 Relevance of the LOO error

As presented in Section 3.3, when the maximal information about y is a set of code evaluations, the error $\|y - \hat{y}^{\text{nest}}\|_{\mathbb{X}}$ can be evaluated by its LOO approximation, ε_{LOO} . In order to reduce the computational cost associated with the evaluation of ε_{LOO} , two alternative estimations of error $\|y - \hat{y}^{\text{nest}}\|_{\mathbb{X}}$, $\hat{\varepsilon}_{\text{LOO}}$ and $\tilde{\varepsilon}_{\text{LOO}}$, have been proposed. In order to underline the relevance of these two LOO errors, Figure 3.4 compares these three errors in the case when $n = 100$ and y is the Ishigami function, whose expression is given by Eq. (3.4.6) (the same kinds of results would have been obtained for other values of n and other expressions of y). In this figure, it can thus be noticed that both approximations $\hat{\varepsilon}_{\text{LOO}}$ and $\tilde{\varepsilon}_{\text{LOO}}$ are very close to $\|y - \hat{y}^{\text{nest}}\|_{\mathbb{X}}$. In general, the approximation $\tilde{\varepsilon}_{\text{LOO}}$ is more conservative, in the sense that

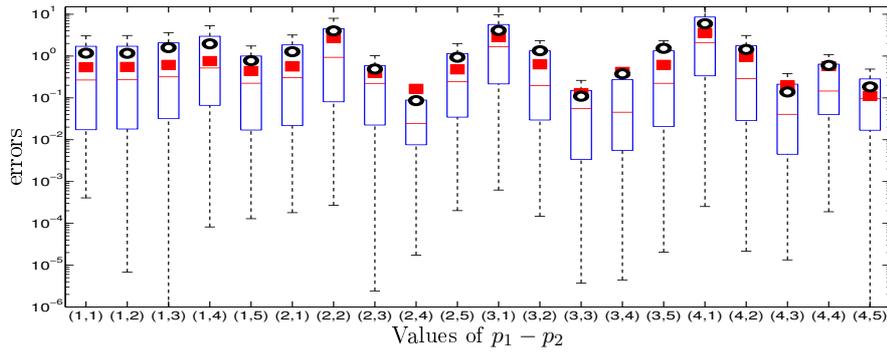
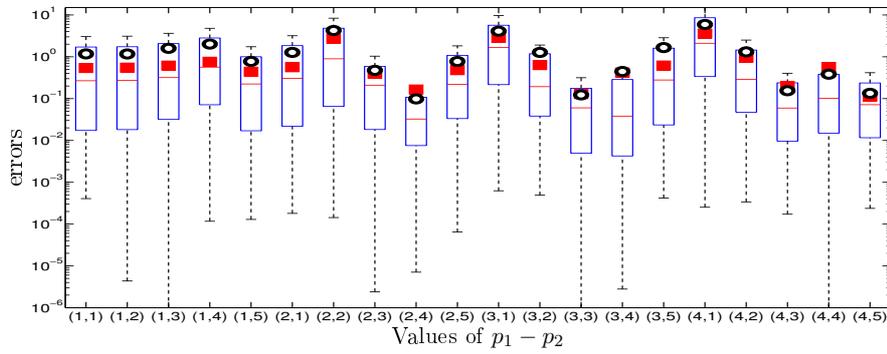

 (a) Case 1: $\hat{\epsilon}_{LOO}$

 (b) Case 2: $\tilde{\epsilon}_{LOO}$

Figure 3.4: Comparisons between error $\|y - \hat{y}^{nest}\|_{\mathbb{X}}$ and its LOO approximations $\hat{\epsilon}_{LOO}$ and $\tilde{\epsilon}_{LOO}$ for the modeling of the Ishigami function from $n = 100$ code evaluations, for $u_2 = d$, $1 \leq p_2 \leq 4$ and $1 \leq p_1 \leq 5$. Red squares: the true values of $\|Y - \hat{y}^{nest}\|_{\mathbb{X}}$. Black circles: the approximated values. In each case, the box-plots correspond to the distributions of $(\hat{e}_n^2, 1 \leq i \leq n)$ and $(\tilde{e}_n^2, 1 \leq i \leq n)$, whose expressions are given by Eqs. (3.3.32) and (3.3.36).

it is less likely that it underestimates $\|y - \hat{y}^{nest}\|_{\mathbb{X}}$. However, as explained in Section 3.3, introducing a linearization around $\hat{\beta}_1$ reduces the risk of being too dependent on $\hat{\beta}_1$, which explains the fact that only small differences can be noticed between $\hat{\epsilon}_{LOO}$ and $\tilde{\epsilon}_{LOO}$.

3.5 Conclusions

One of the main objectives of this part was to propose an alternative parametrization of the polynomial trends for the Gaussian process regression. This parametrization, which is based on the composition of two polynomials, allows us to span high dimensional polynomial spaces with a reduced number of parameters. Hence, it has been shown on a series of examples that this approach can be very useful, especially when confronted to the modeling of complex functions with very little information.

In particular, this approach can allow us to find back (or take into account) a potential nested structure of the code.

However, identifying relevant values for these parameters is not easy. In this chapter, these parameters are identified from a two-steps approach. First, their maximum-likelihood estimates are searched from the resolution of the optimization problem. An iterative algorithm has been proposed to approximate the solutions of this problem. Then, a linearization around these values is carried out, in order to find back the usual formalism of GPR, and to minimize the sensitivity of the results to these values.

In spite of all these adaptations, when the input dimension becomes high ($d > 10$), and when a lot of code evaluations are available ($n > 100d$), it appears that the value of p_1 is often equal to 1. Such a value for p_1 corresponds to the "LAR+UK" configuration, which would mean that, in that case, the nested structure is not necessary. This can be due to the fact that the considered quantity of interest does not present a nested structure, or to the fact that the numerical complexity of the optimization problems associated with the nested representation is too high. Increasing the robustness of the proposed iterative algorithm, as well as proposing more efficient methods to solve the introduced optimization problems are thus possible extensions of the present chapter.

Trying to increase the sparsity of the proposed nested representation could also be a good idea, especially to enable the proposed method to deal with systems with higher values of d . Coupling the proposed nested representation to dedicated penalization techniques seems promising for future work.

Chapter 4

Gaussian process regression of two nested codes with scalar output

In this chapter, we focus on the case of two nested codes with scalar outputs. We now assume that observations of the intermediary variable are available. We therefore consider the following system:

$$\begin{array}{ccc} & \mathbf{x}_2 & \\ & \searrow & \\ \mathbf{x}_1 & \rightarrow & y_1(\mathbf{x}_1) \quad \nearrow \\ & & \end{array} \quad y_{\text{nest}}(\mathbf{x}_{\text{nest}}) := y_2(y_1(\mathbf{x}_1), \mathbf{x}_2), \quad (4.0.1)$$

where \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_{nest} are low dimensional vectors and y_1 , y_2 and y_{nest} are scalars. The work presented in this chapter has been published in Marque-Pucheu et al. [2018]. The framework of the Gaussian process regression is considered (see Chapter 1 for further details). We propose an innovative Gaussian process based sequential design for the case of two nested code with scalar outputs.

4.1 Introduction

Thanks to computing power increase, the certification and the design of complex systems rely more and more on simulation. To this end, predictive codes are needed, which have generally to be evaluated at a large number of input points. When the computational cost of these codes is high, surrogate models are introduced to emulate their responses. A lot of industrial issues involve multi-physics phenomena, which can be associated with a series of computer codes. However, when these code networks are used for optimization, uncertainty quantification, or risk analysis purposes, they are generally considered a single code. In that case, all the inputs characterizing the system of interest are gathered in a single input vector, and little attention is paid to the potential intermediate results. When trying to emulate such code networks, this is clearly sub-optimal, as much information is lost in the statistical learning, so that too many evaluations of each code are likely to be required to get a satisfying prediction precision.

In this chapter, we focus on the case of two nested computer codes, where the output of the first code is one of the inputs of the second code. We assume that these two computer codes are deterministic, but expensive to evaluate. To predict the value of this nested code at an unobserved point, a Bayesian formalism [Robert, 2007] is adopted in the following. Each computer code is *a priori* modeled by a Gaussian process, and the idea is to identify the posterior distribution of the combination of these two processes given a limited number of evaluations of the two codes. The Gaussian process hypothesis is widely used in computer experiments ([Sacks et al., 1989; Santner et al., 2003; Rasmussen and Williams, 2006; Kennedy and O’Hagan, 2000, 2001; Berger et al., 2001; Paulo, 2005; Kleijnen, 2017]), as it allows a very

good trade-off between error control, complexity, and efficiency. The two main issues of this approach, also called Kriging, concern the choice of the statistical properties of the Gaussian processes that are used, and the choice of the points where to evaluate the codes. When a single computer code is considered, several methods exist to add one new point or a batch of new points sequentially to an already existing Design of Experiments. Depending on the purpose, optimization or reconstruction of the objective function on its whole input set, the criteria are based on the mean, variance or covariance of the predictor ([Sacks et al., 1989; Santner et al., 2003; Bect et al., 2012; Echard et al., 2011; Chevalier et al., 2014]). Given that our aim is to predict the output of the nested code on its whole input set, sequential designs based on a reduction of the integrated prediction variance (IMSE) are an appropriate choice. In the case of a single code, the variance expression can be explicitly derived under mild restrictive conditions on the mean and the covariance of the prior Gaussian distribution. The adaptation of these selection criteria to the case of two nested codes is not direct. Indeed, the combination of two Gaussian processes is not Gaussian, so that the prediction variance is much more complicated to estimate. The challenges posed by the composition of two Gaussian processes have been studied in the Deep Gaussian processes literature and the proposed methods are based on the Monte-Carlo computation of the likelihood of the nested Gaussian processes [Perdikaris et al., 2017] or on the computation of a lower bound of this likelihood [Damianou and Lawrence, 2013]. The composition of Gaussian processes can also be used in the multi-fidelity framework [Perdikaris et al., 2017]. This framework enables to use several levels of convergence of a simulator (for example in a finite element model a coarse mesh corresponds to the low fidelity simulator and the finer mesh corresponds to the high-fidelity simulator) and therefore to have a trade-off between accuracy and computation time [Kennedy and O’Hagan, 2000; Le Gratiet, 2013; Le Gratiet and Garnier, 2014; Picheny and Ginsbourger, 2013; Tuo et al., 2014].

Moreover, if the two codes can be launched separately, the selection criterion has also to indicate which one of the two codes to launch. The sequential designs are based on the prediction variance, which has to be computed at a large number of points. To reduce the computational cost associated with these computations, we propose several adaptations of the Gaussian Process formalism to the nested case. These adaptations make it possible to compute the two first statistical moments of the nested code output predictor exactly or quickly. Then, original sequential selection criteria are introduced, which try to exploit as much as possible the nested structure of the studied codes. In particular, these criteria are able to integrate the fact that the computational costs associated with the evaluation of each code can be different.

The outline of this chapter is the following. Section 4.2 presents the theoretical framework of the Gaussian process-based surrogate models, its generalization to the nested case, and introduces two selection criteria based on the prediction variance to reduce the prediction uncertainty sequentially. Section 4.3 introduces a series of simplifications to allow a quick computation of the prediction variance. In Section 4.4, the presented methods are applied to two examples.

The technical proofs of the results presented in the following sections are given in Section 4.6.

4.2 Surrogate modeling for two nested computer codes

4.2.1 General framework

Let \mathcal{S} be a system which is characterized by a vector of input parameters, $\mathbf{x}_{\text{nest}} \in \mathbb{X}_{\text{nest}}$. Let $y_{\text{nest}} : \mathbb{X}_{\text{nest}} \rightarrow \mathbb{R}$ be a deterministic mapping that is used to analyze the studied system. In this chapter, we focus on the case where the function $\mathbf{x}_{\text{nest}} \mapsto y_{\text{nest}}(\mathbf{x}_{\text{nest}})$ can be modeled by

two nested codes. Two quantities of interest, y_1 and y_2 , are thus introduced to characterize these two codes, which are supposed to be two real-valued continuous functions on their respective definition domains \mathbb{X}_1 and $\mathbb{R} \times \mathbb{X}_2$. Given these two functions, the nested code is defined as follows:

$$\begin{array}{ccc} & \mathbf{x}_2 \in \mathbb{X}_2 & \\ & \searrow & \\ \mathbf{x}_1 \in \mathbb{X}_1 & \rightarrow & y_1(\mathbf{x}_1) \in \mathbb{R} \quad \nearrow \\ & & y_{\text{nest}}(\mathbf{x}_{\text{nest}}) := y_2(y_1(\mathbf{x}_1), \mathbf{x}_2) \in \mathbb{R}, \end{array} \quad (4.2.1)$$

where $\mathbf{x}_{\text{nest}} := (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{X}_{\text{nest}} = \mathbb{X}_1 \times \mathbb{X}_2$. The sets \mathbb{X}_1 and \mathbb{X}_2 are moreover supposed to be two compact subsets of \mathbb{R}^{d_1} and \mathbb{R}^{d_2} respectively, where d_1 and d_2 are two positive integers. In theory, the definition domains may be unbounded, but the reduction to compact sets enables the square integrability of y_{nest} on \mathbb{X}_{nest} . If they are unbounded, it is possible to introduce weighted L_2 spaces.

Given a limited number of evaluations of y_1 and y_2 , the objective is to accurately predict y_{nest} on the whole input set.

4.2.2 Gaussian process-based surrogate models

4.2.2.1 Background

The Gaussian process regression (GPR), or Kriging, is a technique that is widely used to replace an expensive computer code by a surrogate model, that is to say a fast to evaluate mathematical function. The GPR is based on the assumption that the two code outputs, y_1 and y_2 , can be seen as the sample paths of two stochastic processes, Y_1 and Y_2 , which are supposed to be Gaussian for the sake of tractability:

$$Y_i(\cdot) \sim \text{GP}(\mu_i(\cdot), C_i(\cdot, \cdot)), \quad i \in \{1, 2\}, \quad (4.2.2)$$

where for all $1 \leq i \leq 2$, μ_i and C_i denote respectively the mean and the covariance functions of Y_i .

Let $\overline{\mathbf{X}}_1^{\text{obs}} := (\overline{\mathbf{x}}_1^{(1)} = \mathbf{x}_1^{(1)}, \dots, \overline{\mathbf{x}}_1^{(n_1)} = \mathbf{x}_1^{(n_1)})$ be a $(n_1 \times d_1)$ -dimensional matrix that gathers n_1 elements of \mathbb{X}_1 and $\overline{\mathbf{X}}_2^{\text{obs}} := (\overline{\mathbf{x}}_2^{(1)} = (\varphi_1^{(1)}, \mathbf{x}_2^{(1)}), \dots, \overline{\mathbf{x}}_2^{(n_2)} = (\varphi_1^{(n_2)}, \mathbf{x}_2^{(n_2)}))$ be a $(n_2 \times d_2)$ -dimensional matrix that gathers n_2 elements of $\mathbb{R} \times \mathbb{X}_2$. Denoting by

$$\mathbf{y}_1^{\text{obs}} := (y_1(\mathbf{x}_1^{(1)}), \dots, y_1(\mathbf{x}_1^{(n_1)})), \text{ and } \mathbf{y}_2^{\text{obs}} := (y_2(\varphi_1^{(1)}, \mathbf{x}_2^{(1)}), \dots, y_2(\varphi_1^{(n_2)}, \mathbf{x}_2^{(n_2)})), \quad (4.2.3)$$

the vectors that gather the evaluations of y_1 and y_2 at these points, it can be shown that:

$$Y_i^c(\cdot) := Y_i(\cdot) \mid \mathbf{y}_i^{\text{obs}} \sim \text{GP}(\mu_i^c(\cdot), C_i^c(\cdot, \cdot)), \quad (4.2.4)$$

and the detailed expressions of the conditioned mean functions, μ_i^c , and the conditioned covariance functions, C_i^c are presented in Eqs. (4.2.11) and (4.2.13) for the "Universal Kriging" framework. For further details on these expressions in other frameworks, the interested reader may refer to Section 1.4.

The relevance of the Gaussian process predictor strongly depends on the definitions of μ_i and C_i . When the only information about y_i is a finite set of evaluations, these functions are generally chosen in general parametric families. In this chapter, functions C_i are chosen in the squared exponential and Matérn-5/2 classes (see Section 1.4 for further details about classical parametric expressions for C_i).

The squared exponential class defines a parametric family of covariance functions that can be written in the form:

$$K_i(\overline{\mathbf{x}}_i, \overline{\mathbf{x}}_i') = \exp\left(-d(\overline{\mathbf{x}}_i, \overline{\mathbf{x}}_i')^2\right), \quad (4.2.5)$$

where:

$$\bar{\mathbf{x}}_i := \begin{cases} \mathbf{x}_1 & \text{if } i = 1, \\ (\varphi_1, \mathbf{x}_2) & \text{if } i = 2, \end{cases} \quad (4.2.6)$$

and $d(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) = \left\| \text{diag}(\boldsymbol{\ell}_i)^{-1} (\bar{\mathbf{x}}_i - \bar{\mathbf{x}}'_i) \right\|$, $\text{diag}(\boldsymbol{\ell}_i)$ denotes a square matrix whose diagonal is equal to the vector $\boldsymbol{\ell}_i$ of correlation lengths.

Regarding the Matérn kernel, we consider the radial Matérn kernel, obtained by substituting the (weighted) Euclidean distance into the 1-dimensional Matérn kernel, and not the tensor product kernel obtained by multiplication of 1-dimensional kernels. So, the covariance functions of the Matérn $\frac{5}{2}$ class can be written in the form:

$$K_i(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) = \left(1 + \sqrt{5}d(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) + \frac{5}{3}d(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i)^2 \right) \exp\left(-\sqrt{5}d(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i)\right). \quad (4.2.7)$$

Linear representations are considered for the mean functions:

$$\mu_i(\cdot) = \mathbf{h}_i(\cdot)^T \boldsymbol{\beta}_i, \quad (4.2.8)$$

where \mathbf{h}_i is a given p_i -dimensional vector of functions (see Chapter 3 for further details on the choice of the basis functions). In the following, the framework of the "Universal Kriging" is adopted, which consists in:

- assuming an (improper) uniform distribution for $\boldsymbol{\beta}_i$,
- conditioning all the results by an estimator of the hyper-parameters that characterize the covariance functions C_i (obtained by cross-validation, as explained below),
- integrating over $\boldsymbol{\beta}_i$ the conditioned distribution of Y_i .

In that case, the distribution of Y_i^c , which is defined by Eq. (4.2.4), is Gaussian, and its statistical moments can explicitly be derived (see Sacks et al. [1989]; Bichon et al. [2008]; Helbert et al. [2009]; Bect et al. [2012]; Perrin et al. [2017]).

If we denote the posterior mean of $\hat{\boldsymbol{\beta}}_i$ by:

$$\hat{\boldsymbol{\beta}}_i := \left[\mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}})^T \right]^{-1} \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} \mathbf{y}_i^{\text{obs}}, \quad (4.2.9)$$

where $\mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}})$ is a $(p_i \times n_i)$ -dimensional matrix, whose j -th column is $\mathbf{h}_i(\bar{\mathbf{x}}_i^{(j)})$, and $C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}})$ is a $(n_i \times n_i)$ -dimensional matrix, such that:

$$\left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)_{jk} = C_i(\bar{\mathbf{x}}_i^{(j)}, \bar{\mathbf{x}}_i^{(k)}), \quad (4.2.10)$$

then the posterior prediction mean and variance can be written:

$$\mu_i^c(\bar{\mathbf{x}}_i) = \mathbf{h}_i(\bar{\mathbf{x}}_i)^T \hat{\boldsymbol{\beta}}_i + C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} \left[\mathbf{y}_i^{\text{obs}} - \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}})^T \hat{\boldsymbol{\beta}}_i \right], \quad (4.2.11)$$

and:

$$(\sigma_i^c(\bar{\mathbf{x}}_i))^2 = C_i^c(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_i), \quad (4.2.12)$$

$$\begin{aligned}
 C_i^c(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) = & C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) - C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{x}}'_i) \\
 & + \left[\mathbf{h}_i(\bar{\mathbf{x}}_i)^T - C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}})^T \right] \\
 & \left[\mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}})^T \right]^{-1} \\
 & \left[\mathbf{h}_i(\bar{\mathbf{x}}'_i) - \mathbf{h}_i(\bar{\mathbf{X}}_i^{\text{obs}}) \left(C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{X}}_i^{\text{obs}}) \right)^{-1} C_i(\bar{\mathbf{X}}_i^{\text{obs}}, \bar{\mathbf{x}}'_i) \right],
 \end{aligned} \tag{4.2.13}$$

where $C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{X}}_i^{\text{obs}})$ is a n_i -dimensional vector and $\left(C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{X}}_i^{\text{obs}}) \right)_k = C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_i^{(k)})$.

In this chapter, the hyperparameters of the covariance functions (see Section 1.4) are estimated for each set of observations by maximizing the Leave-One-Out log predictive probability (see Rasmussen and Williams [2006], Chapter 5, and Bachoc [2013a,b]).

4.2.2.2 Coupling the surrogate models of the two codes

According to Eq. (4.2.1), the nested code, $\mathbf{x}_{\text{nest}} \mapsto y_{\text{nest}}(\mathbf{x}_{\text{nest}})$, can thus be seen as a particular realization of the conditioned process Y_{nest}^c , so that for all $(\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{X}_1 \times \mathbb{X}_2$,

$$Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) := Y_2^c(Y_1^c(\mathbf{x}_1), \mathbf{x}_2). \tag{4.2.14}$$

Under this Gaussian formalism, the best prediction of y_{nest} at any unobserved point $\mathbf{x}_{\text{nest}} = (\mathbf{x}_1, \mathbf{x}_2)$ in $\mathbb{X}_1 \times \mathbb{X}_2$ is given by the mean value of $Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2)$, whereas its variance can be used to characterize the confidence in the prediction. As explained in Section 4.1, there is no reason for Y_{nest}^c to be Gaussian, but according to Proposition 4.2.1, the first- and second-order moments at a given input point can be obtained by computing two one-dimensional integrals with respect to a Gaussian measure.

Proposition 4.2.1. *For all $(\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{X}_1 \times \mathbb{X}_2$, if $\xi \sim \mathcal{N}(0, 1)$, then:*

$$\mathbb{E}[Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2)] = \mathbb{E}[\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1)\xi, \mathbf{x}_2)], \tag{4.2.15}$$

$$\mathbb{E}\left[(Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2))^2\right] = \mathbb{E}\left[\begin{aligned} & \{\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1)\xi, \mathbf{x}_2)\}^2 \\ & + \{\sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1)\xi, \mathbf{x}_2)\}^2 \end{aligned}\right]. \tag{4.2.16}$$

The proof of this Proposition can be found in Section 4.6.

The computation of these moments can be done by quadrature rules or by Monte-Carlo methods ([Baker, 1977]). However, the computation time can be expensive, especially if the moments have to be computed at a large number of points.

Note that the proposed predictor for y_{nest} can be built using observations of y_1 or y_2 alone and not only observations of y_{nest} . It can take into account the partial information. If the two codes can be launched separately, this property will be particularly useful for the sequential enrichment of the initial design of experiments, since the variance of Y_{nest}^c can be reduced by evaluating y_1 or y_2 alone.

4.2.3 Sequential designs for the improvement of Gaussian process predictors

The relevance of the predictor Y_{nest}^c strongly depends on the space filling properties of the sets gathering the inputs of the available observations of y_1 and y_2 , which are generally called Designs of Experiments (DoE). Space-filling Latin hypercube sampling (LHS) or quasi-Monte-Carlo sampling are generally chosen to define such *a priori* DoE ([Fang and Lin, 2003; Fang et al., 2006; Perrin and Cannamela, 2017]). The relevance of the predictor can then be improved by adding new points to an already existing DoE, as the higher the values of n_1 and n_2 , the more chance there is for $\|\mathbb{E}[Y_{\text{nest}}^c] - y_{\text{nest}}\|_{\mathbb{X}_{\text{nest}}}^2$ to be small.

In the case of a single code, the existing selection criteria are based on the prediction variance [Sacks et al., 1989; Santner et al., 2003; Bect et al., 2012; Gramacy and Lian, 2012], the prediction mean [Hu and Ludkovski, 2017] or both [Echard et al., 2011] or the covariance between the observations [Sacks et al., 1989; Santner et al., 2003] and depend on the goal of the experiments: optimization, or reconstruction of the objective function on its whole input domain.

In this chapter the objective is to predict the output of the nested code on its whole input domain. So, a stepwise uncertainty reduction (SUR) [Chevalier et al., 2014] strategy is adopted in order to define criteria to add a new point. The proposed criteria are based on a minimization of the IMSE (integral of the prediction variance over the input domain) or on a maximization of the reduction of IMSE per unit of computational time. Some criteria that enable to take into account the different costs of several computer codes exist, for example in the multi-fidelity framework [Stroh et al., 2017] or multi-objective constraints [Perrin, 2016], but their adaptation to the case of two nested codes is not direct.

The use of IMSE is simplified by some properties of the Gaussian processes. Indeed, if Z is a Gaussian process that is indexed by \mathbf{x} in \mathbb{X} , the variance of the conditioned random variable $Z(\mathbf{x}) | Z(\mathbf{x}^{\text{new}})$, where \mathbf{x} and \mathbf{x}^{new} are any elements of \mathbb{X} , does not depend on the (unknown) value of $Z(\mathbf{x}^{\text{new}})$. So, this variance can be denoted by abuse of notation $\mathbb{V}[Z(\mathbf{x}) | \mathbf{x}^{\text{new}}]$. To minimize the global uncertainty over Z at a reduced computational cost, a natural approach would consist in searching the value of \mathbf{x}^{new} so that

$$\int_{\mathbb{X}} \mathbb{V}[Z(\mathbf{x}) | \mathbf{x}^{\text{new}}] d\mathbf{x} \quad (4.2.17)$$

is minimal (under the condition that this integral exists).

In the nested case, we also have to choose to which code to add a new observation point. To this end, let τ_1 and τ_2 be the numerical costs (in CPU time for instance) that are associated with the evaluations of y_1 and y_2 respectively. For the sake of simplicity, we assume that these numerical costs are independent on the value of the input parameters, and that they are *a priori* known. Two selection criteria are eventually proposed to optimize the relevance of the predictor of the nested code output sequentially. To simplify the reading, the following notation is proposed:

$$(\tilde{\mathbf{x}}_i, \tilde{\mathbb{X}}_i) := \begin{cases} (\mathbf{x}_1^*, \mathbb{X}_1) & \text{if } i = 1, \\ ((\varphi_1^*, \mathbf{x}_2^*), \mu_1^c(\mathbb{X}_1) \times \mathbb{X}_2) & \text{if } i = 2, \\ ((\mathbf{x}_1^*, \mathbf{x}_2^*), \mathbb{X}_1 \times \mathbb{X}_2) & \text{if } i = 3, \end{cases} \quad (4.2.18)$$

where $\mathbf{x}_1^* \in \mathbb{X}_1$, $\varphi_1^* \in \mu_1^c(\mathbb{X}_1)$ and $\mathbf{x}_2^* \in \mathbb{X}_2$ and we denote by $\mathbb{V}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}}) | \tilde{\mathbf{x}}_i)$ the variance of $Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})$ under the hypothesis that the code(s) corresponding to the new point $\tilde{\mathbf{x}}_i$ is (are) evaluated at this point (in practice, we remind that these code evaluations are not required

for the estimation of this variance). This variance can be defined by:

$$\mathbb{V}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})|\tilde{\mathbf{x}}_i) := \begin{cases} \mathbb{V}(Y_2(Y_1(\mathbf{x}_1), \mathbf{x}_2) | \mathbf{y}_1^{\text{obs}}, \mathbf{y}_2^{\text{obs}}, y_i(\tilde{\mathbf{x}}_i)), & i \in \{1, 2\}, \\ \mathbb{V}(Y_2(Y_1(\mathbf{x}_1), \mathbf{x}_2) | \mathbf{y}_1^{\text{obs}}, \mathbf{y}_2^{\text{obs}}, y_{\text{nest}}(\tilde{\mathbf{x}}_i)), & i = 3, \end{cases} \quad (4.2.19)$$

with $\mathbf{x}_{\text{nest}} := (\mathbf{x}_1, \mathbf{x}_2)$.

- First, the chained I-optimal criterion selects the best point in $\mathbb{X}_1 \times \mathbb{X}_2$ to minimize the integrated variance of the predictor of the nested code:

$$\tilde{\mathbf{x}}_3^{\text{new}} = \underset{\tilde{\mathbf{x}}_3 \in \tilde{\mathbb{X}}_3}{\text{argmin}} \int_{\mathbb{X}_{\text{nest}}} \mathbb{V}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})|\tilde{\mathbf{x}}_3) d\mathbf{x}_{\text{nest}}. \quad (4.2.20)$$

Such a criterion is *a priori* adapted to the case where it is not possible to run independently the codes 1 and 2.

- Secondly, the best I-optimal criterion selects the best among the candidates in \mathbb{X}_1 and \mathbb{X}_2 in order to maximize the decrease per unit of computational cost of the integrated prediction variance of the nested code:

$$(\tilde{i}^{\text{new}}, \tilde{\mathbf{x}}_{\tilde{i}^{\text{new}}}^{\text{new}}) = \underset{\tilde{\mathbf{x}}_i \in \tilde{\mathbb{X}}_i, i \in \{1, 2\}}{\text{argmax}} \frac{1}{\tau_i} \times \int_{\mathbb{X}_{\text{nest}}} [\mathbb{V}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})) - \mathbb{V}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})|\tilde{\mathbf{x}}_i)] d\mathbf{x}_{\text{nest}}. \quad (4.2.21)$$

In that case, the difference in the computational costs is taken into account, and a linear expected improvement per unit of computational cost is assumed for the sake of simplicity.

For each new observation of the first code, the hyperparameters of the covariance function C_1 are re-estimated. In the same way, for each new observation of the second code, the hyperparameters of the covariance function C_2 are re-estimated.

An initial set of observations is necessary to estimate the hyperparameters of the covariance functions C_1 and C_2 and therefore to compute the prediction variance and the proposed sequential design criteria. This initial set will be chosen as a maximin LHS design on \mathbb{X}_{nest} .

4.3 Fast computation of the variance of the predictor of the nested code

As explained in Section 4.2.3, choosing the position of the new point requires to compute the value of $\text{Var}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})|\tilde{\mathbf{x}}_i)$ for each potential value of $\tilde{\mathbf{x}}_i$ in $\tilde{\mathbb{X}}_i$ and for a grid or a sample of \mathbf{x}_{nest} used in a quadrature formula or an empirical average to approximate the integral in \mathbf{x}_{nest} of Eqs. (4.2.21) and (4.2.20).

For a given \mathbf{x}_{nest} , the variance is theoretically given by Eqs. (4.2.15) and (4.2.16). If a quadrature rule or a Monte Carlo approach is used to approximate the variance, then the optimization procedure becomes prohibitively expensive from the computational point of view. To circumvent this problem, we present in this section several approaches to make the computation of $\text{Var}(Y_{\text{nest}}^c(\mathbf{x}_{\text{nest}})|\tilde{\mathbf{x}}_i)$ explicit, and therefore extremely fast to compute.

4.3.1 Explicit derivation of the two first statistical moments of the predictor

Lemma 4.3.1. *If $X \sim \mathcal{N}(\mu, \sigma^2)$ and $g(x, a, b, c) := x^a \exp(bx + cx^2)$, $(a, b, c) \in \mathbb{N} \times \mathbb{R}^2$, then, under the condition that $1 - 2c\sigma^2 > 0$, the mean of $g(X, a, b, c)$ can be computed analytically, and its expression is given by Eq. (4.6.1).*

Lemma 4.3.2. *If $g(x, a, b, c) := x^a \exp(bx + cx^2)$, $(a, b, c) \in \mathbb{N} \times \mathbb{R}^2$, then*

$$g(x, a_i, b_i, c_i) g(x, a_j, b_j, c_j) = g(x, a_i + a_j, b_i + b_j, c_i + c_j), \quad (4.3.1)$$

where $(a_i, b_i, c_i) \in \mathbb{N} \times \mathbb{R}^2$ and $(a_j, b_j, c_j) \in \mathbb{N} \times \mathbb{R}^2$.

Proposition 4.3.1. *Using the notations of the Universal Kriging framework that is introduced in Section 4.2.2, if:*

1. for $1 \leq k \leq p_2$ the mean function $(\mathbf{h}_2)_k$ is of the form:

$$(\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_k = m_k(\mathbf{x}_2) \varphi_1^{a_k}, \quad (4.3.2)$$

where m_k is a deterministic function from \mathbb{X}_2 to \mathbb{R} and $a_k \in \mathbb{N}$,

2. the covariance function C_2 is squared exponential, i.e. an element of the squared exponential class,

then the conditional moments of order 1 and 2 of $Y_{nest}^c(\mathbf{x}_1, \mathbf{x}_2)$, which are defined by Eqs. (4.2.15) and (4.2.16) can be calculated analytically using Lemmas 4.3.1 and 4.3.2. Moreover, the expression of the first order moment is given by Eqs. (4.6.5) and (4.6.1) and the one of the second order moment is given by Eqs. (4.6.8) and (4.6.1).

The proof of this Proposition can be found in Section 4.6.

In other words, if the prior of the Gaussian process modeling the function y_2 has a trend which is a polynomial of φ_1 , with coefficients as functions of \mathbf{x}_2 , and a covariance function of the squared exponential class, then the moments of order 1 and 2 of the coupling of the predictors of the two codes can be computed explicitly.

In particular, if the process associated with y_2 has a constant or zero mean and a squared exponential (i.e. Gaussian) covariance, then the mean and the variance of the coupling of the predictors of y_1 and y_2 can be computed analytically.

However, the use of a squared exponential covariance function is based on the assumption of infinite differentiability of the second code. This assumption is not necessarily verified.

Besides, the method cannot be applied to the case of more than two codes. Indeed, in the case of three codes, the coupling of the Gaussian predictors of the two first codes is no longer Gaussian. Even if the Gaussian process modeling the third code has a squared exponential covariance and a polynomial trend with respect to the output of the second code, the analytical method cannot be applied because the predictor of the output of the chain of the two first codes is not Gaussian.

4.3.2 Linearized approach

In the cases where the conditions for Proposition 4.3.1 are not fulfilled (or if more than two codes are considered), another approach is proposed in this section, which is based on a linearization of the process modeling the nested code. Indeed, for $i \in \{1, 2\}$, let ε_i^c be the Gaussian process so that:

$$Y_i^c = \mu_i^c + \varepsilon_i^c. \quad (4.3.3)$$

By construction, ε_i^c is the residual prediction uncertainty once Y_i has been conditioned by n_i evaluations of y_i . We remind that the two Gaussian processes Y_i are statistically independent, so Y_i^c and therefore ε_i^c are statistically independent. Under the condition that n_1 is large enough for Y_1^c being a reliable statistical model for y_1 , then ε_1^c is small.

Proposition 4.3.2. *If:*

1. *the predictor of a nested computer code can be written $Y_{nest}^c(\mathbf{x}_1, \mathbf{x}_2) := Y_2^c(Y_1^c(\mathbf{x}_1), \mathbf{x}_2)$, where Y_i^c are independent Gaussian processes which can be written as $Y_i^c = \mu_i^c + \varepsilon_i^c$, where $\varepsilon_i^c \sim GP(0, C_i^c)$, $i \in \{1, 2\}$,*
2. *and ε_1^c is small enough for the linearization to be valid,*

then the predictor of the nested computer code can be defined as a Gaussian process with the following mean and covariance functions:

$$\begin{aligned} \mu_{nest}^c(\mathbf{x}_1, \mathbf{x}_2) &= \mu_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2), \\ C_{nest}^c((\mathbf{x}_1, \mathbf{x}_2), (\mathbf{x}'_1, \mathbf{x}'_2)) &= C_2^c((\mu_1^c(\mathbf{x}_1), \mathbf{x}_2), (\mu_1^c(\mathbf{x}'_1), \mathbf{x}'_2)) \\ &\quad + \frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}'_1), \mathbf{x}'_2) C_1^c(\mathbf{x}_1, \mathbf{x}'_1), \end{aligned} \quad (4.3.4)$$

where μ_i^c , $i \in 1, 2$ is given by Eq. (4.2.11) and C_i^c , $i \in 1, 2$ is given by Eq. (4.2.13) and $\frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2)$ is given by Eq. (4.6.13).

It can also be written that $Y_{nest}^c = \mu_{nest}^c + \varepsilon_{nest}^c$, with:

$$\varepsilon_{nest}^c(\mathbf{x}_1, \mathbf{x}_2) = \frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \varepsilon_1^c(\mathbf{x}_1) + \varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2). \quad (4.3.5)$$

The proof of this Proposition can be found in Section 4.6.

Corollary 4.3.3. *In the framework of Universal Kriging for Y_1^c and Y_2^c with explicit basis functions \mathbf{h}_i and covariance functions C_i , $i \in \{1, 2\}$, if the derivatives $\frac{\partial \mathbf{h}_2}{\partial \varphi_1}(\varphi_1, \mathbf{x}_2)$ and $\frac{\partial C_2}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{obs})$ can be computed explicitly, then the predictor of the nested computer code can be defined, thanks to a linearization, as a Gaussian process with explicit mean and covariance functions. In particular, if the covariance function C_2 is in the Matérn $\frac{5}{2}$ or squared exponential classes, the derivative $\frac{\partial C_2}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{obs})$ can be computed analytically, and the associated expressions are given in Eqs. (4.6.18) and (4.6.21).*

The proof of this Corollary can be found in Section 4.6.

Corollary 4.3.4. *According to Eqs. (4.3.5), (4.2.21) and (4.2.20), if the predictor of the nested code is obtained with the linearized method, then, thanks to the independence between ε_1^c and ε_2^c , the selection criteria of the sequential designs can be written:*

- for the chained I-optimal design:

$$\begin{aligned}
 (\mathbf{x}_1^{new}, \mathbf{x}_2^{new}) = \underset{(\mathbf{x}_1^*, \mathbf{x}_2^*) \in \mathbb{X}_1 \times \mathbb{X}_2}{\operatorname{argmin}} & \int_{\mathbb{X}_{nest}} \left(\frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \right)^2 \mathbb{V}[\varepsilon_1^c(\mathbf{x}_1) | \mathbf{x}_1^*] d\mathbf{x}_1 d\mathbf{x}_2, \\
 & + \int_{\mathbb{X}_{nest}} \mathbb{V}[\varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) | \mu_1^c(\mathbf{x}_1^*), \mathbf{x}_2^*] d\mathbf{x}_1 d\mathbf{x}_2,
 \end{aligned} \tag{4.3.6}$$

where $\frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2)$ is given by Eq. (4.6.13),

- for the best I-optimal design:

$$(i^{new}, \mathbf{x}_i^{new}) = \underset{\tilde{\mathbf{x}}_i \in \tilde{\mathbb{X}}_i, i \in \{1,2\}}{\operatorname{argmax}} \frac{1}{\tau_i} \mathcal{V}_i(\tilde{\mathbf{x}}_i), \tag{4.3.7}$$

where:

$$\mathcal{V}_1(\tilde{\mathbf{x}}_1) = \int_{\mathbb{X}_{nest}} \left(\frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \right)^2 (\mathbb{V}[\varepsilon_1^c(\mathbf{x}_1)] - \mathbb{V}[\varepsilon_1^c(\mathbf{x}_1) | \tilde{\mathbf{x}}_1]) d\mathbf{x}_1 d\mathbf{x}_2, \tag{4.3.8}$$

$$\mathcal{V}_2(\tilde{\mathbf{x}}_2) = \int_{\mathbb{X}_{nest}} (\mathbb{V}[\varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2)] - \mathbb{V}[\varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) | \tilde{\mathbf{x}}_2]) d\mathbf{x}_1 d\mathbf{x}_2. \tag{4.3.9}$$

The proof of this Corollary can be found in Section 4.6.

Hence, thanks to the proposed linearization, and the fact that the conditional distribution of a Gaussian process is still Gaussian with updated first and second order moments, the variance of $Y_{nest}^c(\mathbf{x}_{nest})$ and the one of $Y_{nest}^c(\mathbf{x}_{nest}) | \tilde{\mathbf{x}}_i$ can be explicitly computed for all $(\mathbf{x}_{nest}, \tilde{\mathbf{x}}_i)$ in $\mathbb{X}_{nest} \times \tilde{\mathbb{X}}_i$. Under the condition that the linearization is valid, this approach can be applied to configurations with more than two nested codes.

However, it can be inferred from equation (4.3.4) that the variance depends on \mathbf{y}_1^{obs} through μ_1^c and \mathbf{y}_2^{obs} through μ_2^c . To circumvent this problem for the computation of the forward variance in the sequential designs, we assume that for a candidate $\tilde{\mathbf{x}}_1$, μ_1^c corresponds to $\mathbb{E}[Y_1 | \mathbf{y}_1^{obs}]$ and by abuse of notation, that $(\sigma_1^c)^2 = C_1^c$ corresponds to $\mathbb{V}[Y_1 | \bar{\mathbf{X}}_1^{obs}, \tilde{\mathbf{x}}_1]$. In the same way, for a candidate $\tilde{\mathbf{x}}_2$, we assume that μ_2^c corresponds to $\mathbb{E}[Y_2 | \mathbf{y}_2^{obs}]$ and by abuse of notation, that $(\sigma_2^c)^2 = C_2^c$ corresponds to $\mathbb{V}[Y_2 | \bar{\mathbf{X}}_2^{obs}, \tilde{\mathbf{x}}_2]$. So, by doing this, we suppose that the estimate of $y_i(\tilde{\mathbf{x}}_i)$ can be replaced by its prediction mean $\mathbb{E}[Y_i(\tilde{\mathbf{x}}_i) | \mathbf{y}_i^{obs}]$, in accordance with the Kriging Believer strategy proposed in Ginsbourger et al. [2010].

4.4 Applications

In this section, the proposed methods are applied to two examples: an analytical one-dimensional one and a multidimensional one.

In particular, the linearized method of Proposition 4.3.2 is compared with the analytical method of Proposition 4.3.1 in terms of prediction accuracy.

The linearized method is compared with the so-called "blind box" method. The blind box method corresponds to the case where the nested computer code is considered as a single computer code. In that case, only the inputs \mathbf{x}_{nest} and the output y_{nest} are taken into account and a Gaussian process regression of this equivalent computer code is done. The intermediary information φ_1 is not taken into account. The Gaussian process Y_{bb} can therefore be defined as follows (see also Perrin et al. [2017]):

$$Y_{bb}(\cdot) \sim GP\left(\mathbf{h}_{bb}(\cdot)^T \boldsymbol{\beta}_{bb}, C_{bb}(\cdot, \cdot)\right), \quad (4.4.1)$$

where

$$\mathbf{h}_{bb}(\mathbf{x}_1, \mathbf{x}_2) = \left(\frac{\partial \mathbf{h}_2}{\partial \varphi_1} \left(\mathbf{h}_1(\mathbf{x}_1)^T \boldsymbol{\beta}_1^*, \mathbf{x}_2 \right)^T \boldsymbol{\beta}_2^* \mathbf{h}_1(\mathbf{x}_1), \mathbf{h}_2 \left(\mathbf{h}_1(\mathbf{x}_1)^T \boldsymbol{\beta}_1^*, \mathbf{x}_2 \right) \right), \quad (4.4.2)$$

$$\boldsymbol{\beta}_{bb} = (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2), \quad (4.4.3)$$

$$(\boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2^*) = \underset{(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)}{\operatorname{argmin}} \sum_{i=1}^n \left[y_2 \left(y_1 \left(\mathbf{x}_1^{(i)} \right), \mathbf{x}_2^{(i)} \right) - \mathbf{h}_2 \left(\mathbf{h}_1 \left(\mathbf{x}_1^{(i)} \right)^T \boldsymbol{\beta}_1, \mathbf{x}_2^{(i)} \right)^T \boldsymbol{\beta}_2 \right]^2, \quad (4.4.4)$$

$n = n_1 = n_2$ and C_{bb} is a stationary covariance function chosen in a parametric family and defined on $\mathbb{X}_{\text{nest}} \times \mathbb{X}_{\text{nest}}$. In order to make the comparison between the blind box and the other methods easier, the mean function is defined as a linearization of the coupling of the mean functions used in the linearized method.

Finally, the performances of the sequential designs are compared with a space filling design (maximin LHS) on \mathbb{X}_{nest} .

4.4.1 Characteristics of the examples

4.4.1.1 Analytical example

In the analytical example, the properties of the mean functions of the Gaussian processes and of the codes are:

$$\mathbf{h}_1(x_1) = \begin{bmatrix} 1 \\ x_1 \\ x_1^2 \end{bmatrix}, \quad \boldsymbol{\beta}_1 = \begin{bmatrix} -2 \\ 0.25 \\ 0.0625 \end{bmatrix}, \quad y_1(x_1) = \mathbf{h}_1(x_1)^T \boldsymbol{\beta}_1 - 0.25 \cos(2\pi x_1), \quad (4.4.5)$$

$$\mathbf{h}_2(\varphi_1) = \begin{bmatrix} 1 \\ \varphi_1 \\ \varphi_1^2 \\ \varphi_1^3 \end{bmatrix}, \quad \boldsymbol{\beta}_2 = \begin{bmatrix} 6 \\ -5 \\ -2 \\ 1 \end{bmatrix}, \quad y_2(\varphi_1) = \mathbf{h}_2(\varphi_1)^T \boldsymbol{\beta}_2 - 0.25 \cos(2\pi \varphi_1), \quad (4.4.6)$$

where $x_1 \in [-7, 7]$. In this example $\mathbb{X}_2 = \emptyset$.

In the analytical example, the covariance functions are squared exponential (i.e. Gaussian). This implies that the Gaussian processes associated with the codes are mean square infinitely differentiable. This enables to apply Proposition 4.3.1 and Proposition 4.3.2 to this example.

4.4.1.2 Hydrodynamic example

In this example, the coupling of two computer codes is considered. The objective is to determine the impact point of a conical projectile.

The first code computes the drag coefficient of a cone divided by the height of the cone. Its inputs are the height and the half-angle of the cone, so the dimension of \mathbf{x}_1 is 2 and $\mathbf{x}_1 \in \left[\frac{\pi}{36}, \frac{\pi}{4} \right] \times [0.2, 2]$.

The second code computes the range of the ballistic trajectory of a cone. Its inputs are the output of the first code, associated with φ_1 , and the initial velocity and angle of the

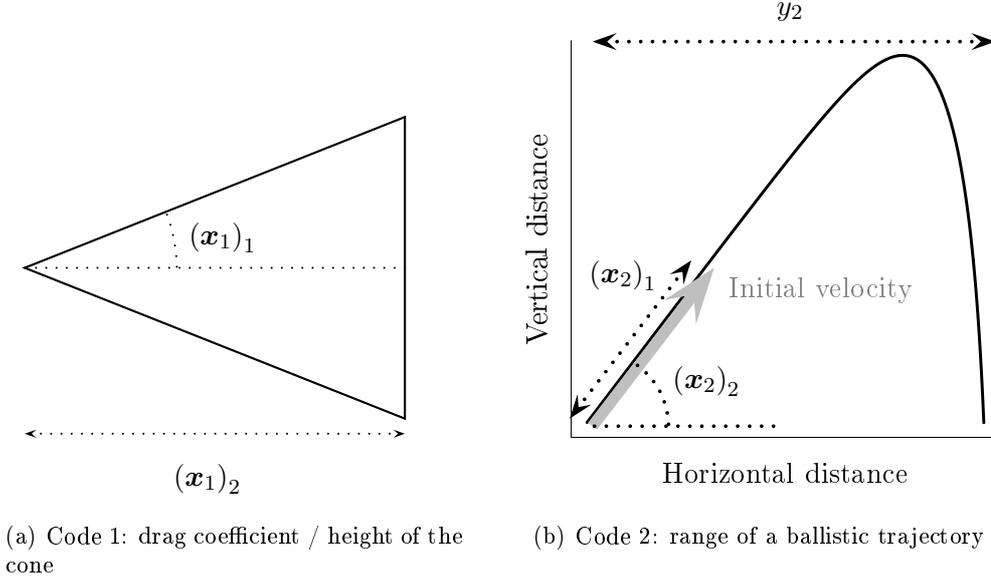


Figure 4.1: Hydrodynamic example: Inputs and outputs of the two codes.

ballistic trajectory of the cone, gathered in \mathbf{x}_2 . The dimension of \mathbf{x}_2 is therefore 2 and $\mathbf{x}_2 \in [1500, 3000] \times \left[\frac{\pi}{12}, \frac{7\pi}{36} \right]$.

Figure 4.1 illustrates the two codes inputs and outputs.

Figure 4.2 presents, for each code, the scatter plots of the variations of the output with respect to the most sensitive components of their inputs. The inputs correspond to a set of 20 points drawn according to a maximin LHS design on \mathbb{X}_{nest} . These figures enable to propose a basis of functions for the prior mean of the processes associated with the two codes.

For the first code, the scatter plots highlight a linear variation with respect to $(\mathbf{x}_1)_1$ and a multiplicative inverse variation with respect to $(\mathbf{x}_1)_2$, so the proposed basis functions are:

$$\mathbf{h}_1(\mathbf{x}_1) = \left(1, (\mathbf{x}_1)_1, \frac{1}{(\mathbf{x}_1)_2} \right)^T. \quad (4.4.7)$$

For the second code, only a multiplicative inverse variation with respect to φ_1 is evident, so the proposed basis functions are:

$$\mathbf{h}_2(\varphi_1, \mathbf{x}_2) = \left(\frac{1}{\max(\varphi_1, \varphi_{1_{\min}})}, 1, 1 \right)^T. \quad (4.4.8)$$

The denominator has a lower bound $\varphi_{1_{\min}}$ in order to avoid any inversion problem around zero. $\varphi_{1_{\min}}$ is set to the small arbitrary value 0.1.

The image plot 4.2(c) represents the UK prediction mean of the first code, obtained with the proposed basis functions. The predicted value of y_1 for the maximum value of $(\mathbf{x}_1)_1$ and the minimum value of $(\mathbf{x}_1)_2$ is high compared with the values of the observations. So, the first code has been evaluated at this input point and gives the value of 3.4, which is consistent with the prediction. This illustrates the relevance of the proposed basis, that is used to extrapolate the prediction at a point with no observations around. The image plot 4.2(e) represents the

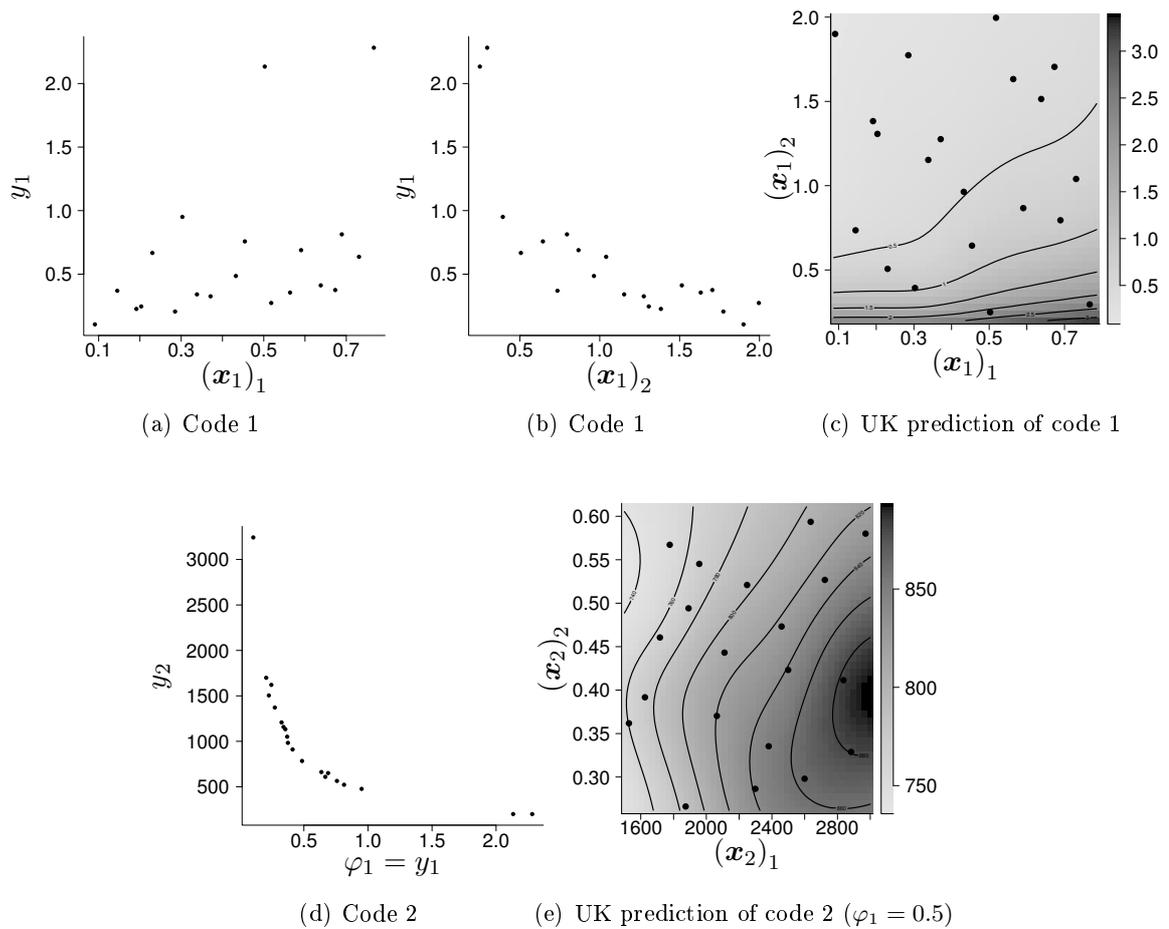


Figure 4.2: Hydrodynamic example: variation of the outputs y_1 and y_2 of the two codes with respect to the most sensitive components of their inputs \mathbf{x}_1 and \mathbf{x}_2 for a set of 20 input points drawn according to a maximin LHS design on \mathbb{X}_{nest} . The image plots present the UK prediction (conditional mean of the GP) of y_1 and y_2 for the same set of observations.

UK prediction mean of the second code, obtained with the proposed basis at a value of 0.5 for φ_1 .

In the hydrodynamic example, the covariance functions are in the Matérn $\frac{5}{2}$ class. This enables to perform the linearization of Proposition 4.3.2 and Corollary 4.3.3.

In both examples, the covariance functions include a non-zero nugget term (see Gramacy and Lee [2012] for further details), that means that they can be written as:

$$C_i(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) = \sigma_i^2 \left[K_i(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i) + g\delta_{\bar{\mathbf{x}}_i = \bar{\mathbf{x}}'_i} \right], \quad (4.4.9)$$

where $\sigma_i \in \mathbb{R}_+$, K_i is chosen in a parametric family (squared exponential or Matérn $\frac{5}{2}$), g is the nugget term whose value is 10^{-6} , and δ is the Kronecker delta function. This non-zero nugget term is used for reasons of numerical stability.

4.4.2 Prediction performance for a given set of observations

A set of validation observations is available. Let $\mathbf{x}_{\text{nest}}^{(1)} \dots \mathbf{x}_{\text{nest}}^{(N_{\text{test}})}$ be N_{test} elements of \mathbb{X}_{nest} . Denoting by $y_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(1)}) \dots y_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(N_{\text{test}})})$ the evaluations of the nested code at these points, the performance criterion of the nested predictor mean, also called error on the mean can be defined as:

$$\text{Error on the mean} = \frac{\sum_{i=1}^{N_{\text{test}}} \left(y_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(i)}) - \hat{y}_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(i)}) \right)^2}{\sum_{i=1}^{N_{\text{test}}} \left(y_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(i)}) - \frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} y_{\text{nest}}(\mathbf{x}_{\text{nest}}^{(j)}) \right)^2}, \quad (4.4.10)$$

where \hat{y}_{nest} denotes a prediction of the nested code, which can be obtained with the analytical, linearized or blind-box method.

For both examples, the validation set of 150 points is drawn according to a maximin LHS on \mathbb{X}_{nest} .

Figure 4.3 presents, for the analytical example, an example of the prediction mean and 95% prediction interval computed with the linearized and the blind box methods. The two predictors are built with the same set of 20 observation points drawn according to a maximin LHS design on \mathbb{X}_{nest} . It can be seen that, with the blind box method, the magnitude of the prediction interval is the same across the input domain and depends only on the distance to the observation points. The prediction interval is too big in the area with small variations and too small in the area with larger variations. On the contrary, taking into account the intermediary observations (with the linearized method here) enables to better take into account the non-stationarity of the variations of the nested code output.

Figure 4.4 presents the error on the mean with the blind box and the linearized methods for both examples, and the analytical method for the analytical example. For all methods, the predictors are built with the same learning sets drawn according to maximin LHS designs on \mathbb{X}_{nest} of increasing size.

The left figure, corresponding to the analytical example, shows the similar accuracies of the prediction means computed with the analytical and linearized methods proposed in Proposition 4.3.1 and Proposition 4.3.2.

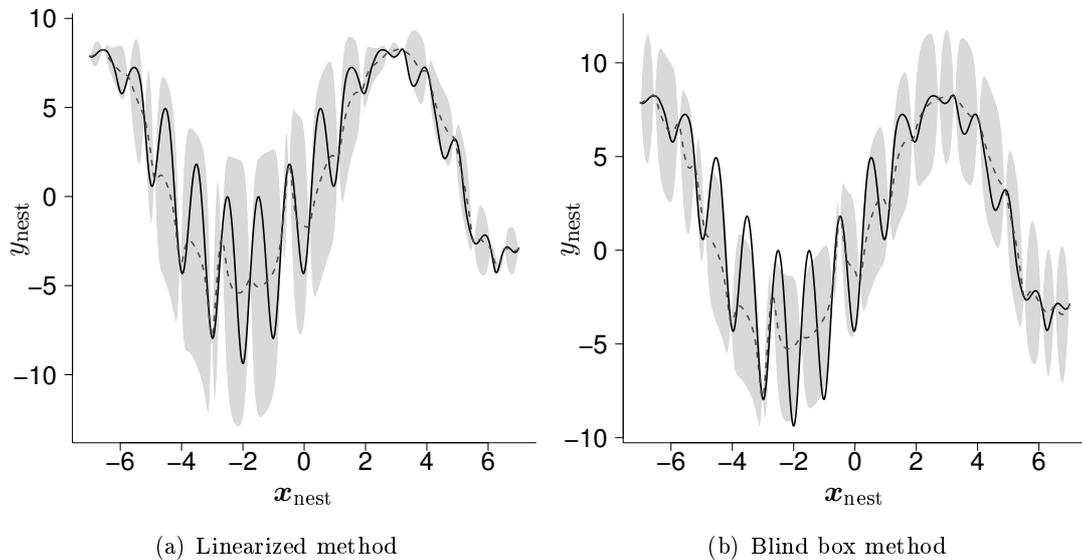


Figure 4.3: Analytical example: Predictors of the nested code obtained with the linearized and the blind box methods. The set of 20 observations is drawn according to a maximin LHS on \mathbb{X}_{nest} . Actual values shown by a continuous line, the prediction mean by a dotted line and the 95% prediction interval by a grey area.

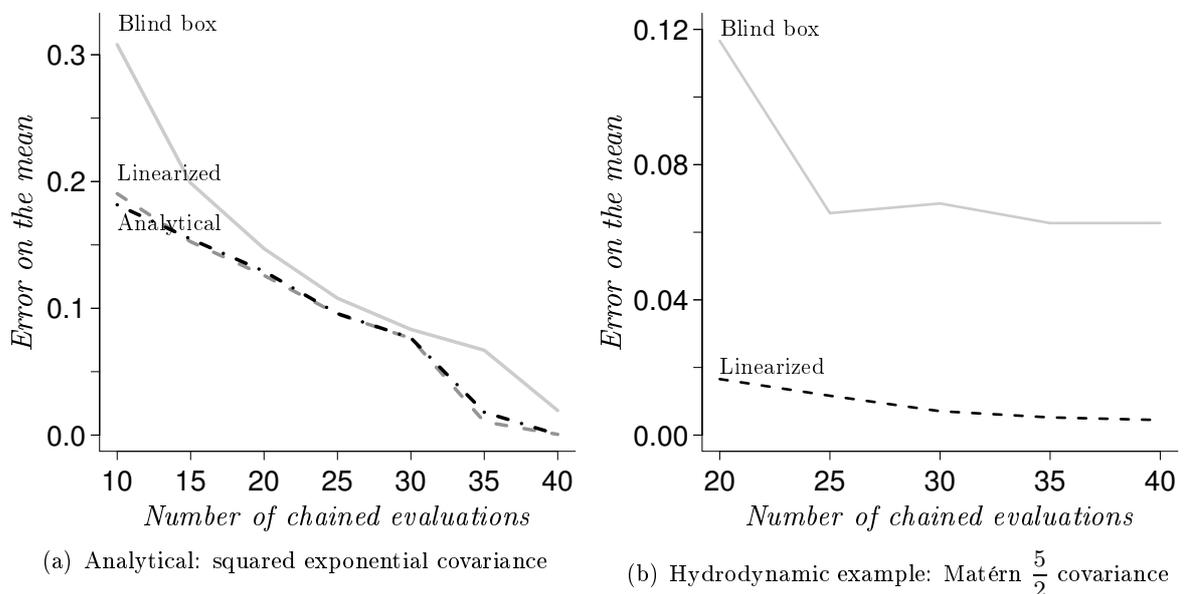


Figure 4.4: Comparison of the prediction mean accuracy for the blind box and the linearized (Proposition 4.3.2) methods, and, in the case of a squared exponential covariance function, the analytical method (Proposition 4.3.1). The curves correspond to the median of 50 draws of maximin LHS designs on $\mathbb{X}_1 \times \mathbb{X}_2$ of increasing size.

For both examples, the precision of the prediction mean is better with the linearized method than with the blind box method, showing the interest of taking into account the intermediary information.

The results show that the analytical and linearized methods lead to the same prediction mean accuracy. As a reminder, the analytical method requires the infinite differentiability of the second code. This assumption is correct for the analytical example but not necessarily for the hydrodynamic example. The linearized method requires the prediction error of the first code to be small enough for the linearization to be valid. Since the prediction error of the first code can be reduced thanks to a sequential enrichment of the initial design, the required assumption of the analytical method is stronger than the one of the linearized method.

Consequently, the linearized method will be used in the remainder of the numerical applications.

4.4.3 Performances of the sequential designs

Figure 4.5 shows an example of the prediction mean and 95% prediction interval of the predictors Y_1^c , Y_2^c and Y_{nest}^c . The predictors Y_1^c and Y_2^c are not built with the same number of observations, so the predictor Y_{nest}^c is built with a different number of observations of the codes 1 and 2. The fact that the number of observations of the two codes can be different will be useful for the sequential designs. Moreover, the estimation of the prediction variance of the nested code is accurate, and that will also be useful for the choice of the new observation point in the sequential designs.

4.4.3.1 With identical computational costs for both codes

Figure 4.6 presents the error on the mean of the linearized predictor for the proposed sequential designs and for maximin LHS designs of increasing size. The initial designs of the sequential strategies are the same maximin LHS designs on \mathbb{X}_{nest} with 10 points for the analytical example and 20 points for the hydrodynamic example. That is why the initial point of the three curves is the same on both line plots. The costs of the two codes are considered to be the same, that is to say $\tau_1 = \tau_2 = 1$. The figure shows the relevance of the proposed sequential designs for improving the prediction mean of the linearized nested predictor, compared with the maximin LHS designs on \mathbb{X}_{nest} .

In the analytical example, the best I-optimal sequential design enables to obtain the most accurate prediction mean at a given computational cost. In the hydrodynamic example, in the first 10 iterations, the best I-optimal design outperforms the chained I-optimal design. After this initial stage, the best I-optimal design calls alternately code 1 and code 2 and becomes equivalent to the chained I-optimal design.

Figure 4.7 shows to which of the two codes the new observations points are added for the best I-optimal sequential design. In both examples, new observation points of the first code are first added.

It seems that the uncertainty propagated from the first code into the second code is predominant at the beginning. The best I-optimal sequential design aims therefore at reducing this uncertainty by first adding new observation points of the first code. Then new observations of both codes are added.

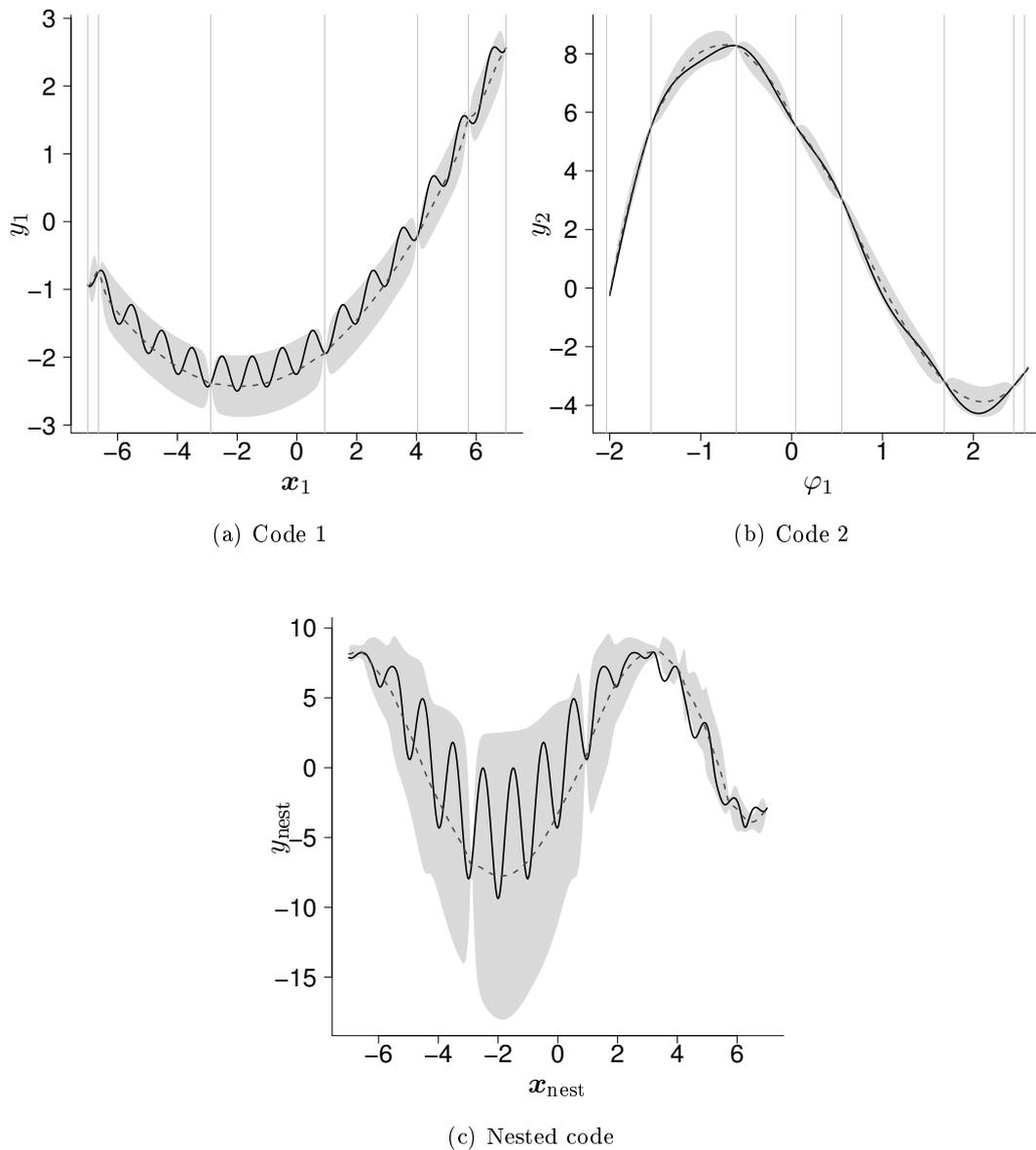


Figure 4.5: Analytical example: an example of the predictors Y_1^c , Y_2^c and Y_{nest}^c . The black line represents the real values of y_1 , y_2 and y_{nest} , the grey area, the 95% prediction interval and the grey dotted line, the prediction mean. The mean and prediction interval of Y_{nest}^c are computed thanks to the linearized method. The vertical lines of the two top plots represent the observations of the two codes, which are drawn according to LHS designs on \mathbb{X}_1 and $\mu_1^c(\mathbb{X}_1)$ of sizes 7 and 8. The number of observations is not the same for both codes.

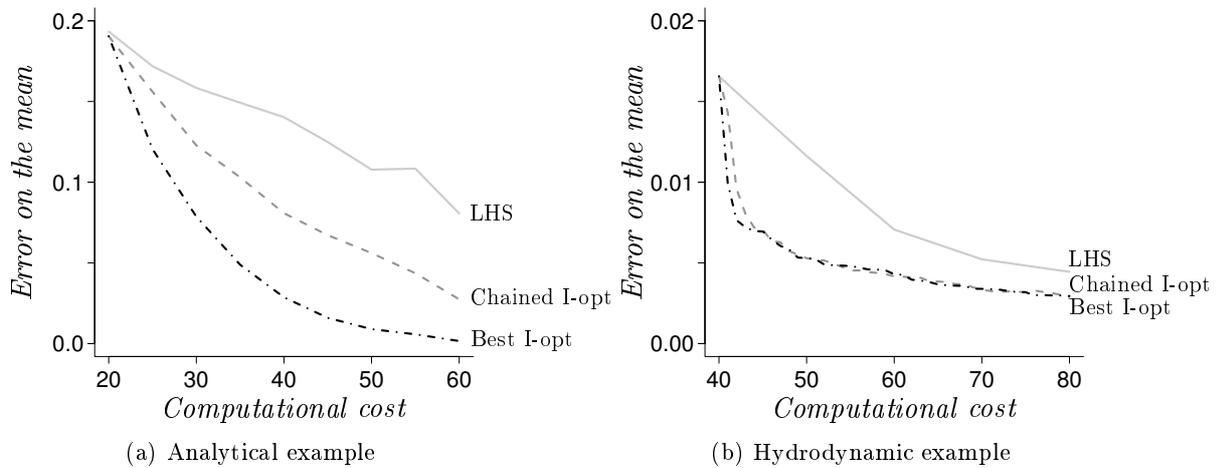


Figure 4.6: Comparison of the prediction mean accuracy of the linearized predictor with the maximin LHS design on \mathbb{X}_{nest} and the sequential designs, for both examples. In the hydrodynamic example, the two curves representing the sequential designs are almost superposed. The initial designs are the same for the three curves, with a size of 10 points for the analytical example and 20 points for the hydrodynamical example. The draw of the maximin LHS design on \mathbb{X}_{nest} is repeated 50 times and the curves present the median of the associated results. The costs of the two codes are assumed to be the same.

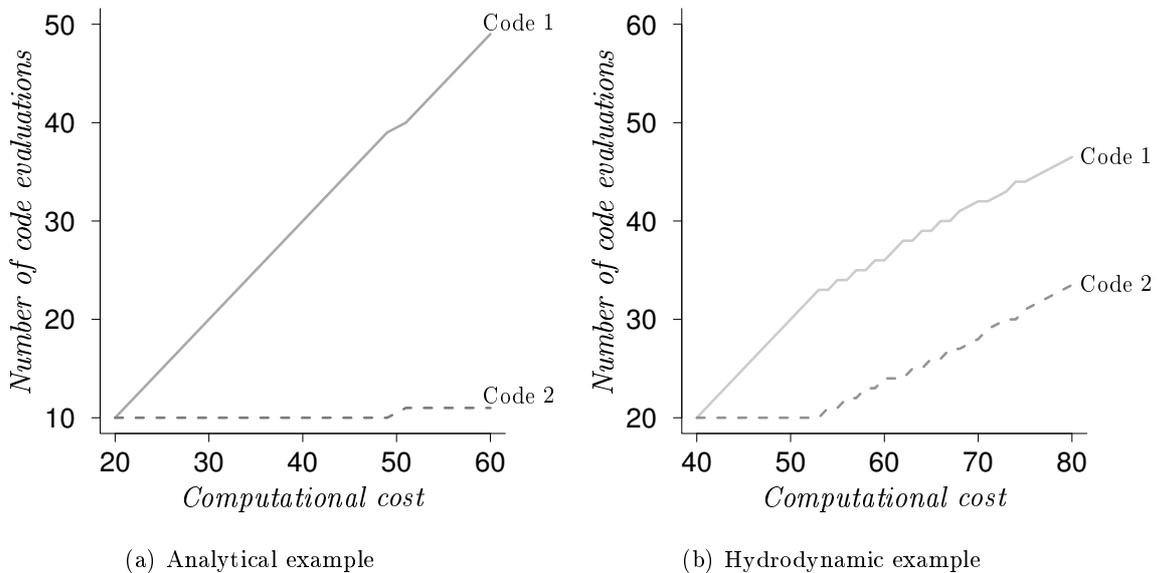


Figure 4.7: Comparison of the number of evaluations of each code in the case of a sequential best I-optimal design applied to both examples. The curves correspond to the median of 50 draws of the initial design. The costs of the two codes are assumed to be the same.

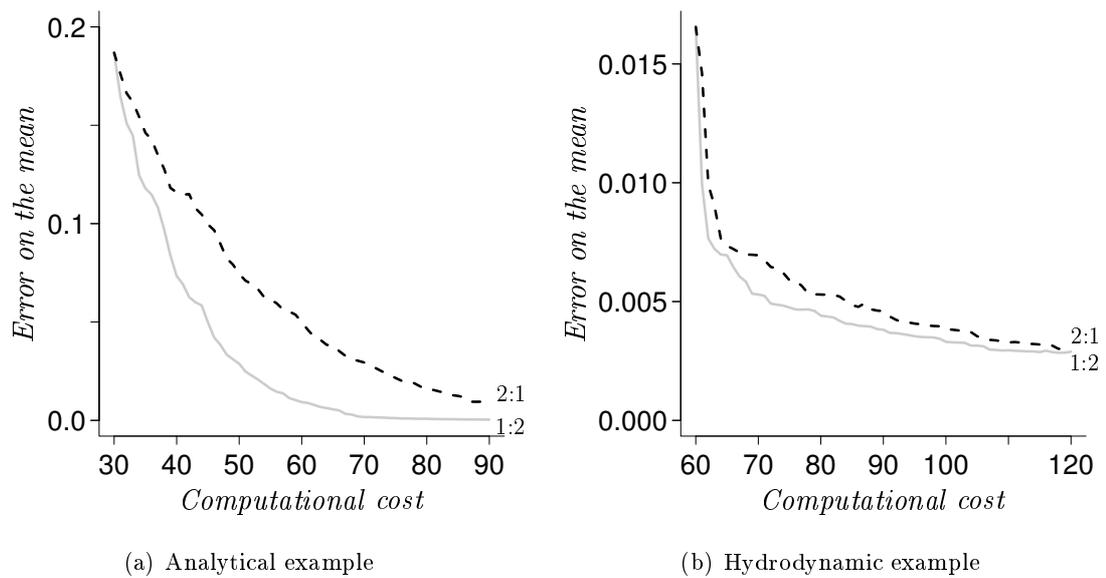


Figure 4.8: Performances of the best I-optimal sequential design in terms of prediction mean accuracy with different computational costs for the two codes. 1:2 \leftrightarrow $\tau_1 = 1$ and $\tau_2 = 2$, 2:1 \leftrightarrow $\tau_1 = 2$ and $\tau_2 = 1$. The curves correspond to the median of 50 draws of the initial maximin LHS design on \mathbb{X}_{nest} . The initial designs are the same for the two curves corresponding to each example and contain 15 observations and 30 observations on both codes for the analytical and the hydrodynamical example.

4.4.3.2 With different computational costs

Figure 4.8 shows the prediction mean accuracy with the best I-optimal sequential design when the costs of the two codes are different. Two cases are presented. The first one corresponds to the case where the cost associated with the first code is twice the one associated with the second code, that is to say $\tau_1 = 2$ and $\tau_2 = 1$, the second one corresponds to the case where the cost associated with the second code is twice the one associated with the first code, that is to say $\tau_1 = 1$ and $\tau_2 = 2$.

It can be seen that for both examples, the prediction accuracy at a given total computational cost is better when the cost of the first code is lower, that is to say when more observation points of the first code can be added for the same computational budget. These results are consistent with those of Figure 4.7.

4.5 Conclusions

In this chapter the formalism of Universal Kriging is adapted to the case of two nested computer codes.

Two methods to compute quickly the mean and variance of the nested code predictor have been proposed. The first one, called "analytical" computes the exact value of the two first moments of the predictor. But it cannot be applied to the coupling of more than two codes. The second one, called "linearized", enables to obtain a Gaussian predictor of the nested code, with mean and variance that can be instantly computed. The approach could be generalized to the coupling of more than two codes.

Both proposed methods take into account the intermediary information, that is to say the output of the first code. A comparison with the reference method, called "blind box", is made. In this method a Gaussian process regression of the block of the two codes is made without considering the intermediary observations. The numerical examples illustrate the interest of taking into account the intermediary information in terms of prediction mean accuracy.

Moreover, two sequential designs are proposed in order to improve the prediction accuracy of the nested predictor. The first one, the "chained" I-optimal sequential design, corresponds to the case where the two codes cannot be launched separately. The second one, the "best" I-optimal sequential design, allows to choose to which of the two codes to add a new observation point and to take into account the different computational costs of the two codes.

The numerical applications show the interest of the sequential designs compared with a space-filling design (maximin LHS). Furthermore, they illustrate the advantage, in terms of prediction mean accuracy, of choosing to which code to add a new observation point compared with simply adding new observation points of the nested code. The results show an amplification of the uncertainties in the chain of codes, leading to the addition of observation points of the first code firstly in the best I-optimal sequential design. It can be assumed that this should be similar with the coupling of more than two codes. In other words, the uncertainty of the beginning of the chain should be reduced as a priority.

4.6 Proofs

4.6.1 Proof of Proposition 4.2.1

According to Eq (4.2.4), one can write:

$$Y_i^c(\mathbf{x}_i) \stackrel{d}{=} \mu_i^c(\mathbf{x}_i) + \sigma_i^c(\mathbf{x}_i) \xi_i, \quad \xi_i \sim \mathcal{N}(0, 1), \quad i \in \{1, 2\},$$

where ξ_1 and ξ_2 are independent according to the independence of the initial processes Y_1 and Y_2 and the fact that $Y_i^c := Y_i | \mathbf{y}_i^{\text{obs}}$.

Therefore, the process modeling the nested code can be written:

$$\begin{aligned} Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) &= Y_2^c(Y_1^c(\mathbf{x}_1), \mathbf{x}_2) \\ &= \mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) + \sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) \xi_2. \end{aligned}$$

Given the independence of ξ_1 and ξ_2 and the fact that $\mathbb{E}(\xi_2) = 0$, it can be inferred that the first moment of Y_{nest}^c can be written:

$$\mathbb{E}(Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2)) = \mathbb{E}(\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2)).$$

By noting that:

$$\begin{aligned} (Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2))^2 &= (Y_2^c(Y_1^c(\mathbf{x}_1), \mathbf{x}_2))^2 \\ &= (\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) + \sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) \xi_2)^2 \\ \bullet &= (\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2))^2 + (\sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2))^2 \xi_2^2 \\ &\quad + 2\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) \sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2) \xi_2, \\ \bullet &\xi_1 \text{ and } \xi_2 \text{ are independent,} \\ \bullet &\mathbb{E}(\xi_2) = 0 \text{ and } \mathbb{E}(\xi_2^2) = 1, \end{aligned}$$

the second moment of Y_{nest}^c can be written:

$$\mathbb{E}\left((Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2))^2\right) = \mathbb{E}\left[\begin{aligned} &(\mu_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2))^2 \\ &+ (\sigma_2^c(\mu_1^c(\mathbf{x}_1) + \sigma_1^c(\mathbf{x}_1) \xi_1, \mathbf{x}_2))^2 \end{aligned}\right].$$

4.6.2 Proof of Lemma 4.3.1

If $X \sim \mathcal{N}(\mu, \sigma^2)$ and $g(x, a, b, c) := x^a \exp[bx + cx^2]$, then the mean of $g(x, a, b, c)$ is equal to:

$$\mathbb{E}[g(X, a, b, c)] = \int_{\mathbb{R}} g(x, a, b, c) \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) dx.$$

It follows that:

$$\begin{aligned}
\mathbb{E}[g(X, a, b, c)] &= \int_{\mathbb{R}} x^a \exp(bx + cx^2) \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) dx \\
&= \exp\left(-\frac{1}{2\sigma^2}\left(\frac{(\sigma^2 b + \mu)^2}{2c\sigma^2 - 1} + \mu^2\right)\right) \\
&\quad \times \int_{\mathbb{R}} x^a \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\frac{1-2c\sigma^2}{\sigma^2}\left(x - \frac{\sigma^2 b + \mu}{1-2c\sigma^2}\right)^2\right) dx \\
&= \exp\left(-\frac{1}{2\sigma^2}\left(\frac{(\sigma^2 b + \mu)^2}{2c\sigma^2 - 1} + \mu^2\right)\right) \frac{1}{\sqrt{1-2c\sigma^2}} \mathbb{E}[X_g^a],
\end{aligned}$$

where $X_g \sim \mathcal{N}\left(\frac{\sigma^2 b + \mu}{1-2c\sigma^2}, \frac{\sigma^2}{1-2c\sigma^2}\right)$, under the condition that $1-2c\sigma^2 > 0$.

Moreover, for $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$, any moment of order k , $k \in \mathbb{N}$, of Y can be computed analytically ([Papoulis and Pillai, 2002]):

$$\mathbb{E}[Y^k] = \sum_{i=0}^{\lfloor \frac{k}{2} \rfloor} \binom{k}{2i} \mu_Y^{k-2i} \frac{(2i)!}{2^i i!} \sigma_Y^{2i}.$$

Hence, given that all the moments of a Gaussian variable can be computed analytically, the mean $\mathbb{E}[g(X, a, b, c)]$ can be computed analytically, and its expression is:

$$\begin{aligned}
\mathbb{E}[g(X, a, b, c)] &= \exp\left(-\frac{1}{2\sigma^2}\left(\frac{(\sigma^2 b + \mu)^2}{2c\sigma^2 - 1} + \mu^2\right)\right) \frac{1}{\sqrt{1-2c\sigma^2}} \\
&\quad \times \sum_{i=0}^{\lfloor \frac{a}{2} \rfloor} \binom{a}{2i} \left(\frac{\sigma^2 b + \mu}{1-2c\sigma^2}\right)^{a-2i} \frac{(2i)!}{2^i i!} \left(\frac{\sigma^2}{1-2c\sigma^2}\right)^i.
\end{aligned} \tag{4.6.1}$$

4.6.3 Proof of Lemma 4.3.2

One has:

$$\begin{aligned}
g(x, a_i, b_i, c_i) g(x, a_j, b_j, c_j) &= x^{a_i} x^{a_j} \exp(b_i x + c_i x^2 + b_j x + c_j x^2) \\
&= x^{a_i + a_j} \exp((b_i + b_j)x + (c_i + c_j)x^2) \\
&= g(x, a_i + a_j, b_i + b_j, c_i + c_j).
\end{aligned}$$

4.6.4 Proof of Proposition 4.3.1

First moment

In the framework of Universal Kriging, according to equation (4.2.11) the conditional mean function of the process modeling the second code can be written:

$$\begin{aligned}\mu_2^c(\varphi_1, \mathbf{x}_2) &= \mathbf{h}_2(\varphi_1, \mathbf{x}_2)^T \widehat{\boldsymbol{\beta}}_2 + C_2((\varphi_1, \mathbf{x}_2), \bar{\mathbf{x}}_2^{\text{obs}}) \mathbf{v}_c \\ &= \sum_{i=1}^{p_2} (\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_i (\widehat{\boldsymbol{\beta}}_2)_i + \sum_{i=1}^{n_2} C_2((\varphi_1, \mathbf{x}_2), (\varphi_1^{(i)}, \mathbf{x}_2^{(i)})) (\mathbf{v}_c)_i \\ &= (1) + (2),\end{aligned}\quad (4.6.2)$$

where $\varphi_1 \sim \mathcal{N}(\mu_1^c, (\sigma_1^c)^2)$, and

$$\mathbf{v}_c = \left(C_2(\bar{\mathbf{X}}_2^{\text{obs}}, \bar{\mathbf{X}}_2^{\text{obs}}) \right)^{-1} \left[\mathbf{y}_2^{\text{obs}} - \mathbf{h}_2(\bar{\mathbf{X}}_2^{\text{obs}})^T \widehat{\boldsymbol{\beta}}_2 \right]. \quad (4.6.3)$$

According to the assumptions of Proposition 4.3.1 the i -th, $i \in \{1, \dots, p_2\}$, component of basis function \mathbf{h}_2 can be written:

$$(\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_i = m_i(\mathbf{x}_2) g(\varphi_1, a_i, 0, 0),$$

with m_i deterministic functions and $g(x, a, b, c) := x^a \exp(bx + cx^2)$, $(a, b, c) \in \mathbb{N} \times \mathbb{R}^2$.

In the same way, the covariance function C_2 is in the squared exponential class, so according to Eq. (4.2.5), it can be written:

$$C_2((\varphi_1, \mathbf{x}_2), (\varphi_1', \mathbf{x}_2')) = \sigma_2^2 k \left(\frac{\varphi_1 - \varphi_1'}{\ell_{\varphi_1}} \right) \prod_{i=1}^{d_2} k \left(\frac{(\mathbf{x}_2)_i - (\mathbf{x}_2')_i}{\ell_i} \right),$$

with $k : x \mapsto \exp(-x^2)$. So, one can write that:

$$\begin{aligned}C_2((\varphi_1, \mathbf{x}_2), (\varphi_1', \mathbf{x}_2')) &= k \left(\frac{\varphi_1 - \varphi_1'}{\ell_{\varphi_1}} \right) \ell(\mathbf{x}_2 - \mathbf{x}_2'), \\ &= \exp \left(- \left(\frac{\varphi_1'}{\ell_{\varphi_1}} \right)^2 \right) g \left(\varphi_1, 0, \frac{2\varphi_1'}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right) \ell(\mathbf{x}_2 - \mathbf{x}_2'),\end{aligned}$$

where ℓ is a deterministic function defined by:

$$\ell(\mathbf{x}_2 - \mathbf{x}_2') = \sigma_2^2 \prod_{i=1}^{d_2} \exp \left(- \left(\frac{(\mathbf{x}_2)_i - (\mathbf{x}_2')_i}{\ell_i} \right)^2 \right), \quad (4.6.4)$$

with $\ell_i, 1 \leq i \leq d_2$ the correlation lengths associated with \mathbf{x}_2 .

So, the terms (1) and (2) of the equation (4.6.2) can be written:

$$\begin{aligned}(1) &= \sum_{i=1}^{p_2} g(\varphi_1, a_i, 0, 0) m_i(\mathbf{x}_2) (\widehat{\boldsymbol{\beta}}_2)_i, \\ (2) &= \sum_{i=1}^{n_2} (\mathbf{v}_c)_i \ell(\mathbf{x}_2 - \mathbf{x}_2^{(i)}) \exp \left(- \left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}} \right)^2 \right) g \left(\varphi_1, 0, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right).\end{aligned}$$

Given that m_i and ℓ are deterministic functions, $\widehat{\boldsymbol{\beta}}_2$, \mathbf{v}_c , $\mathbf{x}_2^{(i)}$ and \mathbf{x}_2 are deterministic vectors, and the $\varphi_1^{(i)}$ are deterministic real numbers, one has:

$$\begin{aligned}\mathbb{E}[(1)] &= \sum_{i=1}^{p_2} \mathbb{E}[g(\varphi_1, a_i, 0, 0)] m_i(\mathbf{x}_2) \left(\widehat{\boldsymbol{\beta}}_2\right)_i, \\ \mathbb{E}[(2)] &= \sum_{i=1}^{n_2} (\mathbf{v}_c)_i \ell\left(\mathbf{x}_2 - \mathbf{x}_2^{(i)}\right) \exp\left(-\left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}}\right)^2\right) \mathbb{E}\left[g\left(\varphi_1, 0, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2}\right)\right].\end{aligned}$$

According to Lemma 4.3.1, and the fact that $1 - 2\left(\frac{-1}{\ell_{\varphi_1}^2}\right)(\sigma_1^c)^2 > 0$, the means $\mathbb{E}[(1)]$ and $\mathbb{E}[(2)]$ can be calculated analytically, and consequently, the mean $\mathbb{E}[\mu_2^c(\varphi_1, \mathbf{x}_2)]$ can be calculated analytically, and its expression is:

$$\begin{aligned}\mathbb{E}[\mu_2^c(\varphi_1, \mathbf{x}_2)] &= \sum_{i=1}^{p_2} \mathbb{E}[g(\varphi_1, a_i, 0, 0)] m_i(\mathbf{x}_2) \left(\widehat{\boldsymbol{\beta}}_2\right)_i \\ &\quad + \sum_{i=1}^{n_2} (\mathbf{v}_c)_i \ell\left(\mathbf{x}_2 - \mathbf{x}_2^{(i)}\right) \exp\left(-\left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}}\right)^2\right) \mathbb{E}\left[g\left(\varphi_1, 0, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2}\right)\right],\end{aligned}\tag{4.6.5}$$

where \mathbf{v}_c is defined by Eq. (4.6.3), $\ell(\mathbf{x}_2 - \mathbf{x}_2')$ is defined by Eq. (4.6.4), ℓ_{φ_1} is the correlation length associated with φ_1 and $\widehat{\boldsymbol{\beta}}_2$ is given by Eq. (4.2.9).

Second moment

From Eq. (4.2.11), (4.2.13) and (4.6.3), one has:

$$\mu_2^c(\varphi_1, \mathbf{x}_2) = \mathbf{h}_2(\varphi_1, \mathbf{x}_2)^T \widehat{\boldsymbol{\beta}}_2 + C_2\left((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}}\right) \mathbf{v}_c,$$

and:

$$\begin{aligned}(\sigma_2^c(\varphi_1, \mathbf{x}_2))^2 &= C_2\left((\varphi_1, \mathbf{x}_2), (\varphi_1, \mathbf{x}_2)\right) - C_2\left((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}}\right) \left(\mathbf{R}_2^{\text{obs}}\right)^{-1} C_2\left(\overline{\mathbf{X}}_2^{\text{obs}}, (\varphi_1, \mathbf{x}_2)\right) + \\ &\quad \left[\mathbf{h}_2(\varphi_1, \mathbf{x}_2)^T - C_2\left((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}}\right) \left(\mathbf{R}_2^{\text{obs}}\right)^{-1} \mathbf{h}_2\left(\overline{\mathbf{X}}_2^{\text{obs}}\right)^T\right] \left[\mathbf{h}_2\left(\overline{\mathbf{X}}_2^{\text{obs}}\right) \left(\mathbf{R}_2^{\text{obs}}\right)^{-1} \mathbf{h}_2\left(\overline{\mathbf{X}}_2^{\text{obs}}\right)^T\right]^{-1} \\ &\quad \left[\mathbf{h}_2(\varphi_1, \mathbf{x}_2) - \mathbf{h}_2\left(\overline{\mathbf{X}}_2^{\text{obs}}\right) \left(\mathbf{R}_2^{\text{obs}}\right)^{-1} C_2\left(\overline{\mathbf{X}}_2^{\text{obs}}, (\varphi_1, \mathbf{x}_2)\right)\right],\end{aligned}$$

where $\mathbf{R}_2^{\text{obs}} = C_2\left(\overline{\mathbf{X}}_2^{\text{obs}}, \overline{\mathbf{X}}_2^{\text{obs}}\right)$.

Hence, it can be written that:

$$\begin{aligned}(\mu_2^c(\varphi_1, \mathbf{x}_2))^2 + (\sigma_2^c(\varphi_1, \mathbf{x}_2))^2 &= \sigma_2^2 + \underbrace{\mathbf{h}_2(\varphi_1, \mathbf{x}_2)^T \mathbf{A}_h \mathbf{h}_2(\varphi_1, \mathbf{x}_2)}_{(1)} \\ &\quad + \underbrace{C_2\left((\varphi_1, \mathbf{x}_2), \overline{\mathbf{x}}_2^{\text{obs}}\right) \mathbf{A}_c C_2\left(\overline{\mathbf{x}}_2^{\text{obs}}, (\varphi_1, \mathbf{x}_2)\right)}_{(2)} \\ &\quad + \underbrace{C_2\left((\varphi_1, \mathbf{x}_2), \overline{\mathbf{x}}_2^{\text{obs}}\right) \mathbf{A}_{ch} \mathbf{h}_2(\varphi_1, \mathbf{x}_2)}_{(3)},\end{aligned}\tag{4.6.6}$$

where:

$$\begin{aligned}
\mathbf{A}_h &= \widehat{\boldsymbol{\beta}}_2 \widehat{\boldsymbol{\beta}}_2^T + \left(\mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right) \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right)^T \right)^{-1}, \\
\mathbf{A}_c &= \mathbf{v}_c \mathbf{v}_c^T - \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} + \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right)^T \left[\mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right) \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right)^T \right]^{-1} \\
&\quad \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right) \left(\mathbf{R}_2^{\text{obs}} \right)^{-1}, \\
\mathbf{A}_{ch} &= 2 \mathbf{v}_c \widehat{\boldsymbol{\beta}}_2^T - 2 \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right)^T \left[\mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right) \left(\mathbf{R}_2^{\text{obs}} \right)^{-1} \mathbf{h}_2 \left(\overline{\mathbf{X}}_2^{\text{obs}} \right)^T \right]^{-1}.
\end{aligned} \tag{4.6.7}$$

According to the assumptions of Proposition 4.3.1 and to lemma 4.3.2, the terms (1), (2) and (3) of the equation (4.6.6) can be rewritten:

$$\begin{aligned}
(1) &= \sum_{i=1}^{p_2} \sum_{j=1}^{p_2} (\mathbf{A}_h)_{ij} (\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_i (\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_j \\
&= \sum_{i=1}^{p_2} \sum_{j=1}^{p_2} (\mathbf{A}_h)_{ij} m_i(\mathbf{x}_2) m_j(\mathbf{x}_2) g(\varphi_1, a_i, 0, 0) g(\varphi_1, a_j, 0, 0) \\
&= \sum_{i=1}^{p_2} \sum_{j=1}^{p_2} (\mathbf{A}_h)_{ij} m_i(\mathbf{x}_2) m_j(\mathbf{x}_2) g(\varphi_1, a_i + a_j, 0, 0), \\
(2) &= \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} (\mathbf{A}_c)_{ij} C_2 \left((\varphi_1, \mathbf{x}_2), (\varphi_1^{(i)}, \mathbf{x}_2^{(i)}) \right) C_2 \left((\varphi_1, \mathbf{x}_2), (\varphi_1^{(j)}, \mathbf{x}_2^{(j)}) \right) \\
&= \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} (\mathbf{A}_c)_{ij} \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(i)} \right) \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(j)} \right) \exp \left(- \frac{(\varphi_1^{(i)})^2 + (\varphi_1^{(j)})^2}{\ell_{\varphi_1}^2} \right) \\
&\quad \times g \left(\varphi_1, 0, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right) g \left(\varphi_1, 0, \frac{2\varphi_1^{(j)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right) \\
&= \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} (\mathbf{A}_c)_{ij} \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(i)} \right) \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(j)} \right) \exp \left(- \frac{(\varphi_1^{(i)})^2 + (\varphi_1^{(j)})^2}{\ell_{\varphi_1}^2} \right) \\
&\quad \times g \left(\varphi_1, 0, 2 \frac{\varphi_1^{(i)} + \varphi_1^{(j)}}{\ell_{\varphi_1}^2}, \frac{-2}{\ell_{\varphi_1}^2} \right), \\
(3) &= \sum_{i=1}^{n_2} \sum_{j=1}^{p_2} (\mathbf{A}_{ch})_{ij} C_2 \left((\varphi_1, \mathbf{x}_2), (\varphi_1^{(i)}, \mathbf{x}_2^{(i)}) \right) (\mathbf{h}_2(\varphi_1, \mathbf{x}_2))_j \\
&= \sum_{i=1}^{n_2} \sum_{j=1}^{p_2} (\mathbf{A}_{ch})_{ij} \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(i)} \right) \exp \left(- \left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}} \right)^2 \right) m_j(\mathbf{x}_2) g \left(\varphi_1, 0, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right) \\
&\quad \times g(\varphi_1, a_j, 0, 0) \\
&= \sum_{i=1}^{n_2} \sum_{j=1}^{p_2} (\mathbf{A}_{ch})_{ij} \ell \left(\mathbf{x}_2 - \mathbf{x}_2^{(i)} \right) \exp \left(- \left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}} \right)^2 \right) m_j(\mathbf{x}_2) g \left(\varphi_1, a_j, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2} \right).
\end{aligned}$$

Given that m_i and ℓ are deterministic functions, \mathbf{x}_2 and $\mathbf{x}_2^{(i)}$ are deterministic vectors, \mathbf{A}_h , \mathbf{A}_c and \mathbf{A}_{ch} deterministic matrices, and $\varphi_1^{(i)}$ and ℓ_{φ_1} are deterministic real numbers, one can write:

$$\mathbb{E} [(1)] = \sum_{i=1}^{p_2} \sum_{j=1}^{p_2} (\mathbf{A}_h)_{ij} m_i(\mathbf{x}_2) m_j(\mathbf{x}_2) \mathbb{E} [g(\varphi_1, a_i + a_j, 0, 0)],$$

$$\mathbb{E} [(2)] = \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} (\mathbf{A}_c)_{ij} \ell(\mathbf{x}_2 - \mathbf{x}_2^{(i)}) \ell(\mathbf{x}_2 - \mathbf{x}_2^{(j)}) \exp\left(-\frac{(\varphi_1^{(i)})^2 + (\varphi_1^{(j)})^2}{\ell_{\varphi_1}^2}\right) \mathbb{E} \left[g\left(\varphi_1, 0, 2\frac{\varphi_1^{(i)} + \varphi_1^{(j)}}{\ell_{\varphi_1}^2}, \frac{-2}{\ell_{\varphi_1}^2}\right) \right],$$

$$\mathbb{E} [(3)] = \sum_{i=1}^{n_2} \sum_{j=1}^{p_2} (\mathbf{A}_{ch})_{ij} \ell(\mathbf{x}_2 - \mathbf{x}_2^{(i)}) \exp\left(-\left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}}\right)^2\right) m_j(\mathbf{x}_2) \mathbb{E} \left[g\left(\varphi_1, a_j, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2}\right) \right].$$

Hence, according to the lemma 4.3.1, the mean $\mathbb{E}[(1)]$ can be computed analytically. In the same way, according to the lemma 4.3.1, and the fact that $1 - 4\left(\frac{-1}{\ell_{\varphi_1}^2}\right)(\sigma_1^c)^2 > 0$ and $1 - 2\left(\frac{-1}{\ell_{\varphi_1}^2}\right)(\sigma_1^c)^2 > 0$, the means $\mathbb{E}[(2)]$ and $\mathbb{E}[(3)]$ can be calculated analytically. Consequently, the mean $\mathbb{E}[(\mu_2^c(\varphi_1, \mathbf{x}_2))^2 + (\sigma_2^c(\varphi_1, \mathbf{x}_2))^2]$ can be calculated analytically, and its expression is:

$$\begin{aligned} \mathbb{E} [(\mu_2^c(\varphi_1, \mathbf{x}_2))^2 + (\sigma_2^c(\varphi_1, \mathbf{x}_2))^2] &= \sigma_2^2 + \sum_{i=1}^{p_2} \sum_{j=1}^{p_2} (\mathbf{A}_h)_{ij} m_i(\mathbf{x}_2) m_j(\mathbf{x}_2) \mathbb{E} [g(\varphi_1, a_i + a_j, 0, 0)] \\ &+ \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} (\mathbf{A}_c)_{ij} \ell(\mathbf{x}_2 - \mathbf{x}_2^{(i)}) \ell(\mathbf{x}_2 - \mathbf{x}_2^{(j)}) \exp\left(-\frac{(\varphi_1^{(i)})^2 + (\varphi_1^{(j)})^2}{\ell_{\varphi_1}^2}\right) \\ &\quad \times \mathbb{E} \left[g\left(\varphi_1, 0, 2\frac{\varphi_1^{(i)} + \varphi_1^{(j)}}{\ell_{\varphi_1}^2}, \frac{-2}{\ell_{\varphi_1}^2}\right) \right] \\ &+ \sum_{i=1}^{n_2} \sum_{j=1}^{p_2} (\mathbf{A}_{ch})_{ij} \ell(\mathbf{x}_2 - \mathbf{x}_2^{(i)}) m_j(\mathbf{x}_2) \exp\left(-\left(\frac{\varphi_1^{(i)}}{\ell_{\varphi_1}}\right)^2\right) \mathbb{E} \left[g\left(\varphi_1, a_j, \frac{2\varphi_1^{(i)}}{\ell_{\varphi_1}^2}, \frac{-1}{\ell_{\varphi_1}^2}\right) \right], \end{aligned} \tag{4.6.8}$$

where \mathbf{A}_h , \mathbf{A}_c and \mathbf{A}_{ch} are defined in Eq. (4.6.7), \mathbf{v}_c is defined in Eq. (4.6.3), $\ell(\mathbf{x}_2 - \mathbf{x}_2')$ is defined by Eq. (4.6.4), ℓ_{φ_1} is the correlation length associated with φ_1 and $\widehat{\beta}_2$ is given by Eq. (4.2.9).

From the two previous paragraphs and Proposition 4.2.1, it can be inferred that, if verifying the assumptions of Proposition 4.3.1, then the first and the second moments of $Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2)$ can be calculated analytically.

4.6.5 Proof of Proposition 4.3.2

If $Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) = Y_2^c(Y_1^c(\mathbf{x}_1), \mathbf{x}_2)$ where $Y_i^c = \mu_i^c + \varepsilon_i^c$, $\varepsilon_i^c \sim \text{GP}(0, C_i^c)$, $i \in \{1, 2\}$, then if ε_1^c is small enough, the process $Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2)$ can be linearized:

$$\begin{aligned} Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) &= \mu_2^c(\mu_1^c(\mathbf{x}_1) + \varepsilon_1^c(\mathbf{x}_1), \mathbf{x}_2) + \varepsilon_2^c(\mu_1^c(\mathbf{x}_1) + \varepsilon_1^c(\mathbf{x}_1), \mathbf{x}_2), \\ &\approx \mu_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) + \frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \varepsilon_1^c(\mathbf{x}_1) + \varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2). \end{aligned}$$

So, one can write:

$$Y_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) \approx \mu_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) + \varepsilon_{\text{nest}}^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2), \quad (4.6.9)$$

with

$$\mu_{\text{nest}}^c(\mathbf{x}_1, \mathbf{x}_2) = \mu_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2), \quad (4.6.10)$$

and

$$\varepsilon_{\text{nest}}^c = \frac{\partial \mu_2^c}{\partial \varphi_1}(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2) \varepsilon_1^c(\mathbf{x}_1) + \varepsilon_2^c(\mu_1^c(\mathbf{x}_1), \mathbf{x}_2). \quad (4.6.11)$$

ε_1^c and ε_2^c are independent centered Gaussian processes, so $\varepsilon_{\text{nest}}^c$ is a centered Gaussian process, whose covariance function, C_{nest}^c , is given by:

$$\begin{aligned} C_{\text{nest}}^c((\mathbf{x}_1, \mathbf{x}_2), (\mathbf{x}'_1, \mathbf{x}'_2)) &= C_2^c((\mu_1^c(\mathbf{x}_1), \mathbf{x}_2), (\mu_1^c(\mathbf{x}'_1), \mathbf{x}'_2)) \\ &\quad + \frac{\partial \mu_2^c}{\partial \varphi_1}((\mu_1^c(\mathbf{x}_1), \mathbf{x}_2)) \frac{\partial \mu_2^c}{\partial \varphi_1}((\mu_1^c(\mathbf{x}'_1), \mathbf{x}'_2)) C_1^c(\mathbf{x}_1, \mathbf{x}'_1). \end{aligned} \quad (4.6.12)$$

From Eqs (4.6.9), (4.6.10), (4.6.11) and (4.6.12), it can be inferred that the predictor of the nested code can be defined as a Gaussian process with mean function μ_{nest}^c defined by Eq. (4.6.10), and covariance function C_{nest}^c defined by Eq. (4.6.12).

Moreover, it follows from Eq. (4.2.11) that:

$$\begin{aligned} \frac{\partial \mu_2^c}{\partial \varphi_1}(\varphi_1, \mathbf{x}_2) &= \left(\frac{\partial \mathbf{h}_2}{\partial \varphi_1}(\varphi_1, \mathbf{x}_2) \right)^T \widehat{\boldsymbol{\beta}}_2 \\ &\quad + \frac{\partial C_2^c}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}}) \left(C_2(\overline{\mathbf{X}}_2^{\text{obs}}, \overline{\mathbf{X}}_2^{\text{obs}}) \right)^{-1} \left[\mathbf{y}_2^{\text{obs}} - \mathbf{h}_2(\overline{\mathbf{X}}_2^{\text{obs}})^T \widehat{\boldsymbol{\beta}}_2 \right]. \end{aligned} \quad (4.6.13)$$

4.6.6 Proof of Corollary 4.3.3

Explicit mean

According to Eq. (4.2.11), if \mathbf{h}_i and C_i can be computed explicitly, then μ_i^c can be computed explicitly. Therefore, according to Eq. (4.3.4), the mean of the Gaussian linearized predictor can be computed explicitly.

Explicit variance

According to Eq. (4.2.13), if \mathbf{h}_i and C_i can be computed explicitly, then C_i^c can be computed explicitly.

According to Eq. (4.6.13), if \mathbf{h}_2 , C_2 and the derivatives $\frac{\partial \mathbf{h}_2}{\partial \varphi_1}(\varphi_1, \mathbf{x}_2)$ and $\frac{\partial C_2^c}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}})$ can be computed explicitly, then the derivative of μ_2^c with respect to φ_1 can be computed explicitly.

Therefore, according to Eq. (4.3.4), the variance of the Gaussian linearized predictor can be computed explicitly.

Hence it follows that, if \mathbf{h}_i and C_i and the derivatives $\frac{\partial \mathbf{h}_2}{\partial \varphi_1}(\varphi_1, \mathbf{x}_2)$ and $\frac{\partial C_2}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}})$ can be computed explicitly, then the mean and the variance of the Gaussian linearized predictor of the nested code can be computed explicitly.

Moreover, the derivative $\frac{\partial C_2}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), \overline{\mathbf{X}}_2^{\text{obs}})$ can be computed explicitly if C_2 is in the squared exponential or the Matérn $\frac{5}{2}$ class, and the associated explicit formulas are given in what follows.

Matérn $\frac{5}{2}$ class

If one denotes by:

$$\begin{aligned} \delta &= d((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2)) \\ &= \sqrt{\frac{(\varphi_1 - \varphi'_1)^2}{\ell_{\varphi_1}^2} + \sum_{i=1}^{d_2} \frac{((\mathbf{x}_2)_i - (\mathbf{x}'_2)_i)^2}{\ell_i^2}}, \end{aligned} \quad (4.6.14)$$

then, according to Eq. (4.2.7), the Matérn kernel can be rewritten:

$$K_{\frac{5}{2}}(\delta) = \left(1 + \sqrt{5}\delta + \frac{5}{3}\delta^2\right) \exp(-\sqrt{5}\delta). \quad (4.6.15)$$

Moreover, one has:

$$\frac{\partial \delta}{\partial \varphi_1} = \frac{\varphi_1 - \varphi'_1}{\ell_{\varphi_1}^2} \frac{1}{\delta}, \quad (4.6.16)$$

and

$$\begin{aligned} \frac{\partial K_{\frac{5}{2}}}{\partial \delta}(\delta) &= \exp(-\sqrt{5}\delta) \left[-\sqrt{5} \left(1 + \sqrt{5}\delta + \frac{5}{3}\delta^2\right) + \sqrt{5} + \frac{10}{3}\delta \right] \\ &= \exp(-\sqrt{5}\delta) \left[-5\delta - \sqrt{5}\frac{5}{3}\delta^2 + \frac{10}{3}\delta \right] \\ &= -\frac{5}{3}\delta (1 + \sqrt{5}\delta) \exp(-\sqrt{5}\delta), \end{aligned} \quad (4.6.17)$$

By noting that in the case of a Matérn $\frac{5}{2}$ kernel:

$$\frac{\partial C_2}{\partial \varphi_1} = \frac{\partial K_{\frac{5}{2}}}{\partial \delta} \frac{\partial \delta}{\partial \varphi_1},$$

the derivative of C_2 with respect to φ_1 is:

$$\begin{aligned} \frac{\partial C_2}{\partial \varphi_1}((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2)) &= -\frac{5}{3} \frac{\varphi_1 - \varphi'_1}{\ell_{\varphi_1}^2} \left[1 + \sqrt{5} d((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2)) \right] \\ &\quad \exp\left[-\sqrt{5} d((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2))\right]. \end{aligned} \quad (4.6.18)$$

Squared exponential class

According to Eq. (4.2.5), the squared exponential kernel can be rewritten:

$$K_{\text{Gauss}}(\delta) = \exp(-\delta^2). \quad (4.6.19)$$

Hence, we have:

$$\frac{\partial K_{\text{Gauss}}}{\partial \delta}(\delta) = -2\delta \exp(-\delta^2). \quad (4.6.20)$$

By noting that, in the case of a squared exponential kernel:

$$\frac{\partial C_2}{\partial \varphi_1} = \frac{\partial K_{\text{Gauss}}}{\partial \delta} \frac{\partial \delta}{\partial \varphi_1},$$

the derivative of C_2 with respect to φ_1 is:

$$\frac{\partial C_2}{\partial \varphi_1} \left((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2) \right) = -2 \frac{\varphi_1 - \varphi'_1}{\ell_{\varphi_1}^2} \exp \left[-d \left((\varphi_1, \mathbf{x}_2), (\varphi'_1, \mathbf{x}'_2) \right)^2 \right]. \quad (4.6.21)$$

Conclusions

This work was motivated by an application case. This application case is the coupling of two computationally expensive codes. The first code is a detonation code and the second code is a structural dynamics code. The two codes have functional (i.e. high dimensional vectorial) outputs. One of the inputs of the second code is the functional output of the first code. The objective was to perform design and certification studies on this system.

The methods used for the design and the certification, like sensitivity analysis, risk analysis or optimization, require a large number of evaluations of the output of the considered system. Considering the high computational cost of the codes, it is in practice impossible to apply these methods directly to the real codes. In this work, we were particularly interested in performing a sensitivity analysis of the output of the nested code with respect to its inputs. The objective of this work was therefore to construct a predictor of the output of the system from a small set of observations, which is accurate on the most likely regions of the input domain of the nested code.

Several difficulties were raised by the surrogate modeling of the considered system:

- There are two codes.
- The two codes are costly, and therefore there will be a few observations available.
- The second code has functional input and output.

This thesis made contributions to help achieve the initially set objective.

The framework of Gaussian process regression, more precisely Universal Kriging in a Bayesian framework, was used.

In a first step, the case of two nested codes with scalar outputs and no intermediary observations was considered. An original parametrization of the mean function of the Gaussian process modeling the nested code was proposed. This parametrization consists of the coupling of two polynomials. It yields a better prediction accuracy than a classical Universal Kriging predictor with a polynomial mean function.

In a second step, the case of two nested codes with scalar outputs and observations of the intermediary variable was considered. A stochastic predictor of the nested code output based on the coupling of Gaussian predictors of the two codes was proposed. The predictor can be constructed from all the types of observations available: those of the nested code, those of the first code and those of the second code. The predictor is non-Gaussian and its mean and variance have to be evaluated with Monte Carlo methods. Furthermore, we proposed two sequential design criteria which aim at improving the accuracy of the predictor on the whole input domain. One of the criteria can take account of the difference between the computational costs of the two codes.

The two sequential design criteria requiring a large number of evaluations of the prediction variance, two adaptations of the predictor were proposed for accelerating the computation of the prediction variance. The first adaptation can lead to closed forms of the mean and the variance of the predictor, if the output is assumed to be infinitely differentiable. The second

one was obtained by proposing a linearization of the coupling of the predictors of the two codes. The predictor of the nested code is then Gaussian with mean and variance functions in closed forms.

In a third and final step, the case of two nested codes with functional outputs and observations of the intermediary variable was considered. An original dimension reduction of the functional input of the second code was proposed. It is based on the approximation of the output of the second code by a linear causal filter and on the projection of the functional input on a basis which is adapted to the linear approximation.

Thanks to this dimension reduction an efficient predictor of the second code is obtained.

Then, similarly to the case of scalar outputs, we proposed a Gaussian predictor of the nested code based on the linearization of the coupling the Gaussian predictors associated with the two codes. Finally, the previously defined sequential design criteria were adapted to the case of codes with functional outputs.

In this thesis, we focused on the surrogate modeling of two nested codes. The study of the surrogate modeling for the coupling of more than two codes or of more complex networks of computer codes is a promising topic. In a non-ringed network, several other relationships between the codes can be found. There can be chains of more than two codes. The output of two different codes can be the inputs of a third code. Besides the case of ringed network could also be studied. Finally, the case of two nested codes with a functional output for the first code and a scalar output for the second code could be studied.

From a practical point of view, the use of parallel computing for the computation of the sequential design criteria with Monte Carlo methods could be useful, especially when the dimension of the input domain is high and the number of Monte Carlo draws too. This could be applied to the case of a one-by-one sequential enrichment of the design. For the case of a batch enrichment, with the addition of $k > 1$ new observations at each step, the number of possible combinations can be very high, which can lead to a high computational burden. The number of possible combinations increases significantly when the number of candidates and k increase. Moreover, the number of candidates is generally higher when the dimension of the inputs is high.

Besides, if we note that the inversion of the covariance matrix of the observations can be expensive when the number of observations is high, it could be interesting to study the possible combination between the proposed linearized predictor and the nested Kriging approach of Rulli re et al. [2018] in order to extend the results obtained to the case of a high number of observations.

The study of optimization strategies for nested codes could also be of great interest. Note that the Expected Improvement criterion presented in Section 1.5 is adapted to the case of a computer code with a scalar output (the quantity to optimize) and its adaptation to the case of a functional output is not direct. If the scalar criterion to be optimized is obtained by a linear transformation of the output, then, thanks to the Gaussianity of the proposed predictor, an enrichment based on the Expected Improvement could be performed.

Bibliography

- P. Abrahamsen. A review of Gaussian random fields and correlation functions. Technical report, Norwegian computing center,, 1997.
- M. Abramowitz and I. Stegun. *Handbook of mathematical functions*. Dover, New York, 1965.
- M. Arnst, R. Ghanem, and C. Soize. Identification of Bayesian posteriors for coefficients of chaos expansions. *Journal of Computational Physics*, 229 (9):3134–3154, 2010.
- F. Bachoc. Cross validation and maximum likelihood estimation of hyper-parameters of Gaussian processes with model misspecification. *Computational Statistics and Data Analysis*, 66:55–69, 2013a.
- F. Bachoc. *Parametric estimation of covariance function in Gaussian-process based Kriging models. Application to uncertainty quantification for computer experiments*. PhD thesis, Université Paris-Diderot - Paris VII, 2013b.
- C. T. H. Baker. *The numerical treatment of integral equations*. Clarendon Press, Oxford, 1977.
- R. A. Bates, R. J. Buck, E. Riccomagno, and H. P. Wynn. Experimental design and observation for large systems. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):77–94, 1996.
- M. J. Bayarri, J. O. Berger, J. Cafeo, G. Garcia-Donato, F. Liu, J. Palomo, R. J. Parthasarathy, R. Paulo, J. Sacks, and D. Walsh. Computer model validation with functional output. *The Annals of Statistics*, 35(5):1874–1906, 2007.
- J. Bect, D. Ginsbourger, L. Li, V. Picheny, and E. Vasquez. Sequential design of computer experiments for the estimation of a probability of failure. *Statistics and Computing*, 22: 773–797, 2012.
- J. O. Berger, V. De Oliveira, and B. Sansó. Objective Bayesian analysis of spatially correlated data. *Journal of the American Statistical Association*, 96(456):1361–1374, 2001.
- B.J. Bichon, M.S. Eldred, L.P. Swiler, S. Mahadevan, and J.M. McFarland. Efficient global reliability analysis for non linear implicit performance functions. *AIAA Journal*, 46(10), 2008.
- I. Bilonis, N. Zabararas, B.A. Konomi, and G. Lin. Multi-output separable Gaussian process: towards an efficient, fully Bayesian paradigm for uncertainty quantification. *Journal of Computational Physics*, 241, 2013.
- G. Blatman and B. Sudret. An adaptive algorithm to build up sparse polynomial chaos expansions for stochastic finite element analysis. *Probabilistic Engineering Mechanics*, 25 (2):183–197, 2010.

- G. Blatman and B. Sudret. Adaptive sparse polynomial chaos expansion based on least angle regression. *Journal of Computational Physics*, 230(6):2345–2367, 2011.
- R. Brent. *Algorithms for Minimization without Derivatives*. Englewood Cliffs N.J.: Prentice-Hall, 1973.
- E. Candes and T. Tao. The dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics*, 35(6):2313–2351, 12 2007.
- C. Chevalier, J. Bect, D. Ginsbourger, and E. Vazquez. Fast parallel Kriging-based stepwise uncertainty reduction with application to the identification of an excursion set. *Technometrics*, 56(4):455–465, 2014. doi: 10.1080/00401706.2013.860918>.
- P.G. Constantine, E. Dow, and Q-Q. Wang. Active Subspace methods in theory and practice: Applications to Kriging surfaces. *SIAM Journal on Scientific Computing*, 36(4):1500–1524, 2014.
- S. Conti, J.P. Gosling, J.E. Oakley, and A. O’Hagan. Gaussian process emulation of dynamic computer codes. *Biometrika*, 96(3):663–676, 2009.
- R.D. Cook and C.J. Nachtsheim. A comparison of algorithms for constructing exact D-optimal designs. *Technometrics*, 22:315–324, 1980.
- A. Damianou and N. D. Lawrence. Deep Gaussian processes. In C. Carvalho and P. Ravikumar, editors, *Proceedings of the Sixteenth International Workshop on Artificial Intelligence and Statistics (AISTATS)*, AISTATS ’13, pages 207–215. JMLR W&CP 31, 2013.
- S. Das, R. Ghanem, and S. Finette. Polynomial chaos representation of spatio-temporal random field from experimental measurements. *J. Comput. Phys.*, 228:8726–8751, 2009.
- O. Dubrule. Cross validation of Kriging in a unique neighborhood. *Mathematical Geology*, 15(6):687–699, 1983.
- B. Echard, N. Gayton, and M. Lemaire. AK-MCS: An active learning reliability method combining Kriging and Monte Carlo simulation. *Structural Safety*, 33:145–154, 2011.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Ann. Stat.*, 32:407–499, 2004.
- M. Efronson. *Mathematical models for digital computers*, volume 1, chapter Multiple regression analysis, pages 191—203. Wiley, 1960.
- G. Elfving. Optimum allocation in linear regression theory. *The Annals of Mathematical Statistics*, 23(2):255–262, 1952.
- K.T. Fang and D.K. Lin. Uniform experimental designs and their applications in industry. *Handbook of Statistics*, 22:131–178, 2003.
- K.T. Fang, R. Li, and A. Sudjianto. *Design and modeling for computer experiments*. Chapman & Hall, Computer Science and Data Analysis Series, London, 2006.
- V.V. Fedorov. *Theory Of Optimal Experiments*. Academic Press, New York, 1972.
- V.V. Fedorov and P. Hackl. *Model-Oriented Design of Experiments*. Springer-Verlag, New York, 1997.
- R. Fisher. *Statistical Methods for Research Workers*. Oliver & Boyd, 1925.

- G.M. Furnival and R.W. Wilson. Regressions by leaps and bounds. *Technometrics*, 16(4): 499–511, 1974.
- D. Geman and B. Jedynak. An active testing model for tracking roads in satellite images. 18: 1 – 14, 02 1996.
- R. Ghanem and P. D. Spanos. Polynomial chaos in stochastic finite elements. *Journal of Applied Mechanics*, 57(1):197–202, 1990.
- R. Ghanem and P. D. Spanos. *Stochastic Finite Elements: A Spectral Approach, rev. ed.* Dover Publications, New York, 2003.
- D.M. Ghiocel and R.G. Ghanem. Stochastic finite-element analysis of seismic soil-structure interaction. *Journal of Engineering Mechanics*, 128(1):66–77, 2002.
- D. Ginsbourger, R. Le Riche, and L. Carraro. *Computational intelligence in expensive optimization problems*, volume 2 of *Adaptation Learning and Optimization*, chapter Kriging is well-suited to parallelize optimization, pages 131–162. Springer Berlin Heidelberg, 2010.
- R. Gramacy and H. Lian. Gaussian process single-index models as emulators for computer experiments. *Technometrics*, 54:1:30–41, 2012.
- R. B. Gramacy and H. K. H. Lee. Cases for the nugget in modeling computer experiments. *Statistics and Computing*, 22:713–722, 2012.
- J. H. Halton. Algorithm 247: Radical-inverse quasi-random point sequence. *Commun. ACM*, 7(12):701–702, 1964.
- J. Hammersley. *Monte Carlo Methods*. Springer Netherlands, 1964.
- M.S. Handcock and M.L. Stein. A Bayesian analysis of Kriging. *Technometrics*, 35:403–4010, 1993.
- T. Hastie and R. Tibshirani. *Generalized additive models*. Chapman & Hall, 1990.
- T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: Data mining, inference and prediction*. Springer, New York, 2001.
- T. Hastie, R. Tibshirani, and Friedman. *Elements of Statistical Learning*. Springer, New York, 2002.
- T. Hastie, J. Taylor, R. Tibshirani, and G. Walther. Forward stagewise regression and the monotone lasso. *Electronic Journal of Statistics*, 1:1–29, 2007.
- C. Helbert, D. Dupuy, and L. Carraro. Assessment of uncertainty in computer experiments, from Universal to Bayesian Kriging. *Applied Stochastic Models in Business and Industry*, 25:99–113, 2009.
- T. Hesterberg, N. H. Choi, L. Meier, and C. Fraley. Least angle and ℓ_1 penalized regression: A review. *Statistics Surveys*, 2:61–93, 2008.
- D. Higdon, J. Gattiker, B. Williams, and M. Rightley. Computer model calibration using high-dimensional output. *Journal of the American Statistical Association*, 103(482):570–583, 2008.
- W. Hoeffding. A class of statistics with asymptotically normal distributions. *The Annals of Mathematical Statistics*, 19:293–325, 1948.

- A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- T. Homma and A. Saltelli. Importance measures in global sensitivity analysis of nonlinear models. *Reliability Engineering & System Safety*, 52(1):1–17, 1996.
- A. Höskuldsson. PLS regression methods. *Journal of chemometrics*, 2(3):211–228, 1988.
- R. Hu and M. Ludkovski. Sequential design for ranking response surfaces. *SIAM/ASA Journal on Uncertainty Quantification*, 5:212–239, 2017.
- J. D. Jakeman, M. S. Eldred, and K. Sargsyan. Enhancing ℓ_1 -minimization estimates of polynomial chaos expansions using basis selection. *Journal of Computational Physics*, 289:18 – 34, 2015.
- G. James, P. Radchenko, and J. Lv. DASSO: Connections between the dantzig selector and lasso. *J. Royal Stat. Soc., Series B*, 71(1):127–142, 2008.
- M. E. Johnson, L. M. Moore, and D. Ylvisaker. Minimax and maximin distance designs. *Journal of Statistical Planning and Inference*, 26(2):131–148, 1990.
- D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box functions. *Biometrika*, 13:455–492, 1998.
- V. R. Joseph, Y. Hung, and A. Sudjianto. Blind Kriging: A new method for developing metamodels. *Journal of mechanical design*, 130(3):031102, 2008.
- M. C. Kennedy and A. O’Hagan. Predicting the output from a complex computer code when fast approximations are available. *Biometrika*, 87:1–13, 2000.
- M. C. Kennedy and A. O’Hagan. Bayesian calibration of computer models. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 63(3):425–464, 2001.
- P. Kersaudy, B. Sudret, N. Varsier, and O. Picon. A new surrogate modeling technique combining Kriging and polynomial chaos expansions - application to uncertainty analysis in computational dosimetry. *Journal of Computational Physics*, 286:103–117, 2015.
- J. Kiefer. Optimum designs in regression problems, ii. *The Annals of Mathematical Statistics*, 32(1):298–325, 1961.
- J. Kiefer and J. Wolfowitz. Optimum designs in regression problems. *The Annals of Mathematical Statistics*, 30(2):271–294, 1959.
- J.P.C. Kleijnen. Regression and Kriging metamodels with their experimental designs in simulation: A review. *European Journal of Operational Research*, 256:1–16, 2017.
- K. Konakli and B. Sudret. Polynomial meta-models with canonical low-rank approximations: Numerical insights and comparison to sparse polynomial chaos expansions. *Journal of Computational Physics*, 321:1144 – 1169, 2016.
- H.J. Kushner. A new method of locating the maximal point of an arbitrary multipeak curve in the presence of noise. *J. Basic Eng.*, 86:97–106, 1964.
- L Le Gratiet. *Multi-fidelity Gaussian process regression for computer experiments*. PhD thesis, Université Paris-Diderot - Paris VII, 2013.

- L. Le Gratiet and J. Garnier. Recursive co-Kriging model for design of computer experiments with multiple levels of fidelity. *International Journal for Uncertainty Quantification*, 4(5):365–386, 2014.
- O.P. Le Maître and O.M. Knio. *Spectral Methods for Uncertainty Quantification*. Springer, 2010.
- R. Lebrun and A. Dutfoy. Do Rosenblatt and Nataf isoprobabilistic transformations really differ? *Probabilistic Engineering Mechanics*, 24(4):577–584, 2009.
- M. Loève. *Probability theory*. Springer, 1955.
- S. Marque-Pucheu, G. Perrin, and J. Garnier. Efficient sequential experimental design for surrogate modeling of nested codes. *ESAIM Probability and Statistics*, 2018.
- G. Matheron and F. Blondel. *Traité de géostatistique appliquée*. Paris, 1962.
- M.D. McKay, R.J. Beckman, and W.J. Conover. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.
- N. Meinshausen, G. Rocha, and B. Yu. A tale of three cousins: Lasso, ℓ_2 -boosting, and dantzig. *Annals Statistics*, 35:2373–2384, 2007.
- Y. Meyer and D. H. Salinger. *Wavelets and operators*, volume 1. Cambridge university press, Cambridge, 1995.
- R. G. Miller. The jackknife - a review. *Biometrika*, 61:1–15, 1974.
- I. Molchanov and S. Zuyev. Steepest descent algorithms in a space of measures. *Statistics and Computing*, 12(2):115–123, 2002.
- S. Nanty, C. Helbert, A. Marrel, N. Pérot, and C. Prieur. Uncertainty quantification for functional dependent random variables. *Computational Statistics*, 32(2):559–583, 2017.
- A. Nataf. Détermination des distributions de probabilité dont les marges sont données. *Comptes Rendus de l'Académie des Sciences*, 225:42–43, 1962.
- H. Niederreiter. Quasi-Monte Carlo methods and pseudo-random numbers. *Bulletin of the American Mathematical Society*, 84(6):957–1041, November 1978.
- A. Nouy. Proper generalized decomposition and separated representations for the numerical solution of high dimensional stochastic problems. *Archives of computational methods in engineering*, 17:403–434, 2010.
- J. Oakley and A. O'Hagan. Bayesian inference for the uncertainty distribution of computer model outputs. *Biometrika*, 89(4):769–784, 2002.
- J. E. Oakley and A. O'Hagan. Probabilistic sensitivity analysis of complex models: A Bayesian approach. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 66(3):751–769, 2004.
- A. O'Hagan. Curvefitting and optimal design for prediction. *J. R. Stat. Soc., Ser. B, Methodol.*, 40(1):1–42, 1978.
- A. Papoulis and S. U. Pillai. *Probability, Random Variables and Stochastic Processes*. McGraw-Hill, Boston, 2002.

- E. Parzen. An approach to time series analysis. *Ann. Math. Stat.*, 32:951–989, 1962.
- R. Paulo. Default priors for Gaussian processes. *Annals of Statistics*, 33(2):556–582, 2005.
- K. Pearson. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 2(11):559–572, 1901.
- P. Perdikaris, M. Raissi, A. Damianou, N. D. Lawrence, and G. E. Karniadakis. Nonlinear information fusion algorithms for data-efficient multi-fidelity modelling. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 473(2198), 2017.
- G. Perrin. Active learning surrogate models for the conception of systems with multiple failure modes. *Reliability Engineering and System Safety*, 149:130–136, 2016.
- G. Perrin. Adaptive calibration of a computer code with time-series output. *Journal of Statistical Planning and Inference*, 2018.
- G. Perrin and C. Cannamela. A repulsion-based method for the definition and the enrichment of optimized space filling designs in constrained input spaces. *Journal de la Société Française de Statistique*, 158(1):37–67, 2017.
- G. Perrin, C. Soize, D. Duhamel, and C. Funfschilling. Identification of polynomial chaos representations in high dimension from a set of realizations. *SIAM Journal on Scientific Computing*, 34(6):2917–2945, 2012.
- G. Perrin, C. Soize, D. Duhamel, and C. Funfschilling. A posteriori error and optimal reduced basis for stochastic processes defined by a finite set of realizations. *SIAM/ASA J. Uncertainty Quantification*, 2:745–762, 2014.
- G. Perrin, C. Soize, S. Marque-Pucheu, and J. Garnier. Nested polynomial trends for the improvement of Gaussian process-based predictors. *Journal of Computational Physics*, 346:389–402, 2017.
- V Picheny and D Ginsbourger. A nonstationary space-time Gaussian process model for partially converged simulations. *SIAM/ASA Journal on Uncertainty Quantification*, 1(1):37–67, 2013.
- V. Picheny, D. Ginsbourger, O. Roustant, R.T Haftka, and N-H. Kim. Adaptive designs of experiments for accurate approximation of a target region. *Journal of Mechanical Design*, 132(7):071008–071008–9, 2010.
- A. Pinkus. *Ridge functions*. Cambridge University Press, Cambridge, 2015.
- C. E. Rasmussen and C. K.I. Williams. *Gaussian processes for machine learning*. The MIT Press, Cambridge, 2006.
- C. Robert. *The Bayesian Choice*. Springer-Verlag New York, New York, 2007.
- M. Rosenblatt. Remarks on a multivariate transformation. *The Annals of Mathematical Statistics*, 23(3):470–472, 1952.
- J. Rougier. Efficient emulators for multivariate deterministic functions. *Journal of Computational and Graphical Statistics*, 17(4):827–843, 2008.
- R. T. Rubinstein and D.P. Kroese. *Simulation and the Monte Carlo method*. John Wiley and Sons, Inc., Hoboken, New Jersey, 2008.

- D. Rullière, N. Durrande, F. Bachoc, and C. Chevalier. Nested Kriging predictions for datasets with a large number of observations. *Statistics and Computing*, 28(4), 2018.
- T.M. Russi. *Uncertainty Quantification with Experimental Data and Complex System Models*. PhD thesis, UC Berkeley, 2010.
- J. Sacks, W.J. Welch, T.J. Mitchell, and H.P. Wynn. Design and analysis of computer experiments. *Statistical Science*, 4(4):409–423, 1989.
- A. Saltelli and I. Sobol. About the use of rank transformation in sensitivity of model output. *Reliability Engineering & System Safety*, 50(3):225 – 239, 1995.
- A. Saltelli, K. Chan, and E. M. Scott. *Sensitivity analysis*. Wiley, New Jersey, 2000.
- T.J. Santner, B.J. Williams, and W.I. Notz. *The Design and Analysis of Computer Experiments*. Springer-Verlag New York, 2003.
- I.M. Sobol. Distribution of points in a cube and approximate evaluation of integrals. *U.S.S.R Comput. Maths. Math. Phys.*, 7:86–112, 1967.
- I.M. Sobol. Sensitivity estimates for nonlinear mathematical models. *Mathematical Modeling & Computational Experiment*, 1:407–414, 1993.
- C. Soize and R. Ghanem. Physical systems with random uncertainties: Chaos representations with arbitrary probability measure. *SIAM Journal on Scientific Computing*, 26:395–410, 2004.
- M.L. Stein. *Interpolation of spatial data: some theory for Kriging*. Springer, New York, 1999.
- R. Stroh, S. Demeyer, N. Fischer, J. Bect, and E. Vazquez. Sequential design of experiments to estimate a probability of exceeding a threshold in a multi-fidelity stochastic simulator. In *61th World Statistics Congress of the International Statistical Institute (ISI 2017)*, Marrakech, Morocco, July 2017.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B*, 58(1):267–288, 1989.
- A. N. Tikhonov and V. Y. Arsenin. *Solution of Ill-posed Problems*. Winston & Sons, Washington, 1977.
- R. Tuo, C.F. Jeff Wu, and D. Yu. Surrogate modeling of computer experiments with different mesh densities. *Technometrics*, 56(3):372–380, 2014.
- J.G. Van der Corput. Verteilungsfunktionen. I. Mitt. *Proc. Akad. Wet. Amsterdam*, 38: 813–821, 1935.
- E. Vazquez and J. Bect. A sequential Bayesian algorithm to estimate a probability of failure. In *15th IFAC Symposium on System Identification, SYSID 2009*, Saint-Malo, France, July 2009.
- B. Williams, D. Higdon, J. Gattiker, L. Moore, M. McKay, and S. Keller-McNulty. Combining experimental data and computer simulations, with an application to flyer plate experiments. *Bayesian Analysis*, 1(4):765–792, 2006.
- H. Wold. *Estimation of Principal Components and Related Models by Iterative Least squares*. Academic Press, 1966.

- C.-F. Wu and H. P. Wynn. The convergence of general step-length algorithms for regular optimum design criteria. *The Annals of Statistics*, 6(6):1273–1285, 11 1978.
- D. Xiu. Fast numerical methods for stochastic computations: a review. *Communications in computational physics*, 5(2-4):242–272, 2009.
- O. Zahm, P. Constantine, C. Prieur, and Y. Marzouk. Gradient-based dimension reduction of multivariate vector-valued functions. *HAL preprint : hal-1801.07922*.