



HAL
open science

Predictive quality of meta-models constructed on the reproducing kernel Hilbert spaces and sensitivity analysis of complex models.

Halaleh Kamari

► **To cite this version:**

Halaleh Kamari. Predictive quality of meta-models constructed on the reproducing kernel Hilbert spaces and sensitivity analysis of complex models.. General Mathematics [math.GM]. Université Paris-Saclay, 2020. English. NNT : 2020UPASE010 . tel-02997897

HAL Id: tel-02997897

<https://theses.hal.science/tel-02997897v1>

Submitted on 10 Nov 2020

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.

Qualité prédictive des méta-modèles construits sur des espaces de Hilbert à noyau auto-reproduisant et analyse de sensibilité des modèles complexes.

Thèse de doctorat de l'Université Paris-Saclay

École doctorale de Mathématique Hadamard (EDMH) n°
574

Spécialité de doctorat: Mathématiques appliquées
Unité de recherche: Université Paris-Saclay, CNRS, Univ Evry,
Laboratoire de Mathématiques et Modélisation d'Evry, 91037,
Evry-Courcouronnes, France.
Réfèrent: : Université d'Evry Val d'Essonne

Thèse présentée et soutenue à Evry, le 06/07/2020, par

Halaleh Kamari

Composition du jury:

Agathe Guilloux Professeur, Université d'Evry-Val-d'Essonne (LaMME)	Présidente
Clémentine Prieur Professeur, Université Grenoble Alpes (AIRSEA)	Rapportrice
Olivier Roustant Professeur, Institut de Mathématiques de Toulouse (INSA)	Rapporteur
Béatrice Laurent-Bonneau Professeur, Institut de Mathématiques de Toulouse (INSA)	Examinatrice
Marie-Luce Taupin Professeur, Université d'Evry-Val-d'Essonne (LaMME)	Directrice
Sylvie Huet Chercheuse (retraitee), INRA (MaIAGE)	Coencadrante

Acknowledgments

Before supervising my thesis, Marie-Luce Taupin was my professor when I was in Master 2. It is actually thanks to her that I met Sylvie Huet and started under their supervision an internship in the unité MaIAGE of INRA in Jouy en Josas. Working with them and along with the rest of the MaIAGE members, inspired me to pursue my studies in this direction, and I am grateful to them for allowing me to do so. I would like to thank my supervisors Sylvie and Marie-Luce for their support, their reactivity and their trust. I have learned a lot from them, both scientifically and in terms of human relations in the world of research.

When I started my internship in MaIAGE, I was really welcomed and I felt as well all along my three first years of thesis. I am grateful that I had the opportunity to discover the members of the unité LaMME in Evry almost in the last year of my thesis. A great thank you to all the members of unité MaIAGE and LaMME. I would like to also thank all the PhD students of MaIAGE and LaMME for having so kindly integrated me into their little teams.

Lastly, I would like to thank my husband Rebaz, my mother, my father, my sisters and my friends who always believed on me, offered me their unconditional love and have supported me all these years.

Contents

1	Introduction	1
1.1	Cadre de travail	1
1.1.1	Introduction sur l'analyse de sensibilité	2
1.1.2	Analyse de sensibilité globale: méthodes basées sur la décomposition de la variance	4
1.1.3	Méta-modélisation	6
1.1.4	Méta-modèles basés sur des espaces à noyaux auto-reproduisants (RKHS)	8
1.1.5	Méthode d'estimation	10
1.2	Résumé du chapitre 3	14
1.2.1	Objectifs et résultats	14
1.2.2	Présentation du modèle	15
1.2.3	Les principaux résultats	16
1.2.4	Travaux antérieurs	19
1.2.5	Outils techniques pour les preuves	20
1.3	Résumé du chapitre 4	21
1.3.1	Objectifs et résultats	21
1.3.2	Présentation du modèle	22
1.3.3	Critère à minimiser	22
1.3.4	Algorithmes	24
1.4	Résumé et perspectives	27
1.4.1	Variables d'entrée non-indépendantes	28
1.4.2	Généralisation au modèle de régression avec erreur log-concave	29
2	Introduction in english	31
2.1	Framework	31
2.1.1	Introduction to the sensitivity analysis	32
2.1.2	Global sensitivity analysis: variance-based methods	34
2.1.3	Meta-modelling	36
2.1.4	Meta-models based on the reproducing kernel Hilbert spaces (RKHS)	37
2.1.5	Estimation method	40
2.2	Summary of Chapter 3	44
2.2.1	Objectives and results	44
2.2.2	Presentation of the model	45
2.2.3	Main results	45
2.2.4	Related works	49
2.2.5	Technical tools for the proofs	50
2.3	Summary of Chapter 4	51
2.3.1	Objectives and results	51
2.3.2	Presentation of the model	52

2.3.3	Criterion to minimize	52
2.3.4	Algorithms	54
2.4	Summary and perspectives	57
2.4.1	Non-independent input variables	57
2.4.2	Generalization to the regression framework with log-concave error	59
3	Risk upper bounds for RKHS ridge group sparse estimator in the regression model with non-Gaussian and non-bounded error	61
3.1	Introduction	61
3.2	Meta-modelling and the RKHS ridge group sparse estimator	65
3.2.1	RKHS construction	65
3.2.2	Approximating the Hoeffding decomposition of m	66
3.2.3	Ridge group sparse procedure and associated estimator	67
3.3	Risk upper bounds	67
3.3.1	Rate of convergence	71
3.4	Main arguments of the proof of Theorem 3.3.1 and motivation for the choice π_α	72
3.4.1	Sketch of the proof	73
3.4.2	Sudakov minoration	77
3.4.3	Concentration inequality	80
3.5	Proof of Theorem 3.3.1	84
3.5.1	Intermediate Lemmas	88
3.5.2	Proof of lemma 3.5.1 to 3.5.4	89
3.5.3	Proofs of intermediate Lemmas	101
3.6	Proof of Corollary 3.3.1	107
Appendix	108
3.A	Proofs of Section 3.4.2	108
3.A.1	Proof of Remark 3.4.1	108
3.A.2	Proof of Corollary 3.4.1	108
3.B	Proofs of Section 3.4.3	109
3.B.1	Proof of Lemma 3.4.2	109
3.B.2	Proof of Remark 3.4.3	110
3.B.3	Proof of Corollary 3.4.2	111
4	Estimate the Hoeffding decomposition of a complex model by solving RKHS ridge group sparse optimization problem	113
4.1	Introduction	113
4.2	Estimation method	119
4.2.1	RKHS ridge group sparse and RKHS group lasso procedures	119
4.2.2	RKHS construction	121
4.2.3	Choice of the tuning parameters	123
4.2.4	Estimation of the Sobol indices	124
4.3	Algorithms	124

4.3.1	Calculation of the Gram matrices	125
4.3.2	Optimization algorithms	127
4.4	Overview of the RKHSMetaMod functions	130
4.4.1	Main RKHSMetaMod functions	131
4.4.2	Companion functions	134
4.5	RKHSMetaMod through examples	137
4.6	Summary and discussion	147
Appendix	149
4.A	More technical details	149
4.A.1	RKHS group lasso algorithm	149
4.A.2	RKHS ridge group sparse algorithm	152
Appendix		154
A	Package ‘RKHSMetaMod’	155
A.1	calc_Kv function	157
A.2	grplasso_q function	159
A.3	mu_max function	161
A.4	pen_MetMod function	163
A.5	PredErr function	166
A.6	RKHSgrplasso function	167
A.7	RKHSMetMod function	169
A.8	RKHSMetMod_qmax function	172
A.9	SI_emp function	175
Bibliography		179

List of Tables

4.1	List of the reproducing kernels used to construct the RKHS \mathcal{H}	125
4.2	List of the input arguments of the <code>RKHSMetMod</code> function.	132
4.3	List of the arguments of the output "Meta-Model" of <code>RKHSMetMod</code> function.	133
4.4	Example 4.5.1: The columns of the table correspond to the different datasets with $n \in \{50, 100, 200\}$ and $d = 5$. Each line of the table, from up to down, gives the value of GPE obtained for each dataset associated with the "matern", "brownian" and "gaussian" kernels, respectively.	139
4.5	Example 4.5.1: The columns of the table correspond to the different datasets with $n \in \{50, 100, 200\}$ and $d = 5$. Each line of the table, from up to down, gives the value of MSE obtained for each dataset associated with the "matern", "brownian" and "gaussian" kernels, respectively.	140
4.6	Example 4.5.1: The first line of the table gives the true values of the Sobol indices $\times 100$. The second line gives the mean of the estimated empirical Sobol indices $\times 100$ greater than 10^{-2} calculated over fifty simulations for $n = 200$ and "matern" kernel. The sum of the Sobol indices is displayed in the last column.	140
4.7	Example 4.5.2: The true values of the Sobol indices $\times 100$ when $d = 10$.	140
4.8	Example 4.5.3: Obtained prediction errors in step 1.	143
4.9	Example 4.5.3: Obtained prediction errors in step 2.	144
4.10	Example 4.5.3: The estimated empirical Sobol indices $\times 100$ greater than 10^{-2} . The last two columns show $\sum_v \hat{S}_v$ and RE, respectively.	144
4.11	Example 4.5.4: The kernel used is "matern". The execution time for the functions <code>RKHSgrplasso</code> and <code>pen_MetMod</code> is displayed in each row for two pair of values of tuning parameters ($\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7)$, $\gamma = 0.01$) on up, and ($\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8)$, $\gamma = 0.01$) on below. In the column $ S_{\hat{f}} $, the number of the active groups associated with each estimated RKHS meta-model is displayed.	145
4.12	Example 4.5.4: Obtained prediction errors.	146
4.13	Example 4.5.4: The estimated empirical Sobol indices $\times 100$ greater than 10^{-2} associated with each estimated RKHS meta-model is printed. The last two columns show $\sum_v \hat{S}_v$ and RE, respectively. We have $\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7)$, $\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8)$ and $\gamma = 0.01$	146

List of Figures

- 4.1 Example 4.5.4: Timing plot for $d = 10$, $n \in \{100, 300, 500, 1000, 2000, 5000\}$, and different functions of the **RKHSMetaMod** package. The execution time for the functions **RKHSgrplasso** and **pen_MetMod** is displayed for two pair of values of tuning parameters ($\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7), \gamma = 0.01$) in solid lines, and ($\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8), \gamma = 0.01$) in dashed lines. 145
- 4.2 On the left, the RKHS meta-model versus the g-function is plotted. On the right, the empirical Sobol indices in the y axis and vMax= 175 groups in the x axis are displayed. 147

Introduction

1.1 Cadre de travail

Considérons un phénomène décrit par un modèle m dépendant de d variables d'entrée $X = (X_1, \dots, X_d)$. Ce modèle m de \mathbb{R}^d vers \mathbb{R} peut être complexe, présenter de fortes non-linéarités et des effets d'interaction d'ordre élevé. Dans le cadre classique de l'analyse de sensibilité, le modèle m peut être calculé en un nombre fini de points.

Lorsque les coordonnées de X sont indépendantes, le modèle m peut se décomposer selon la décomposition dite de Hoeffding. Quand la loi des coordonnées de X est connue, la décomposition de Hoeffding de m permet d'effectuer l'analyse de sensibilité, et plus précisément de calculer les indices de Sobol de m (Sobol (2001), Saltelli et al. (2009)). Cependant, le calcul de ces indices peut être très difficile, voire impossible, surtout lorsque le nombre de variables d'entrée d est grand (Iooss (2011)).

Une approche récente consiste à approcher m par un méta-modèle additif impliquant les coordonnées de X ainsi que leurs interactions, comme proposée par Durrande et al. (2013). Ce méta-modèle, noté f^* , est la projection orthogonale de m sur un espace de Hilbert à noyau auto-reproduisant (RKHS), noté \mathcal{H} . L'espace \mathcal{H} est associé à un noyau dite d'ANOVA qui est défini de façon à obtenir l'expression analytique des termes de la décomposition de Hoeffding des fonctions de \mathcal{H} . Comme f^* est la projection orthogonale de m sur \mathcal{H} , chaque terme de sa décomposition est une approximation du terme associé de la décomposition de Hoeffding de m .

Lorsque le nombre de variables d'entrée d est grand, le nombre total de termes dans la décomposition de Hoeffding de f^* devient très élevé. Une solution consiste à calculer une approximation sparse ou parcimonieuse de f^* en utilisant le critère des moindres carrés pénalisé comme dans le modèle de régression non-paramétrique. Sparse ou parcimonieuse au sens où le nombre de termes non-nuls dans la décomposition de Hoeffding de f^* est contrôlé.

Dans cette thèse, deux cadres sont considérés: l'analyse de sensibilité où $m(X)$ est calculable en tout point X , et le modèle de régression où m est inconnu et ne peut donc pas être calculé.

Dans le second cas, pour un X donné, $m(X)$ est observable à une erreur près. Ainsi, on dispose de l'observation Y telle que,

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0. \quad (1.1)$$

Comme dans le cadre classique de l'analyse de sensibilité, l'idée est d'approcher la décomposition de Hoeffding de m par le méta-modèle f^* , puis de calculer un

estimateur sparse de f^* en utilisant des méthodes d'estimation non-paramétriques. Cet estimateur, noté \hat{f} , est la solution d'un problème de minimisation des moindres carrés pénalisé par une fonction de pénalité qui favorise sparsité et régularité. La construction de l'estimateur permet d'estimer facilement les indices de Sobol de m .

Cette thèse se compose d'une partie théorique et d'une partie pratique:

- Dans la partie théorique, j'ai établi les majorations du risque empirique L^2 et du risque quadratique de l'estimateur \hat{f} d'un modèle de régression (voir l'équation (1.1)) où l'erreur ε est non-gaussienne et non-bornée. Il s'agit des bornes supérieures par rapport à la norme empirique L^2 et à la norme L^2 pour la distance entre la fonction réelle m et son estimation \hat{f} dans le RKHS \mathcal{H} . Cette partie est présentée dans le chapitre 3.
- Dans la partie pratique, j'ai développé un package R appelé **RKHSMetaMod**, pour la mise en œuvre des méthodes d'estimation du méta-modèle f^* d'un modèle m . Ce package s'applique indifféremment dans le cas où le modèle m est calculable et le cas du modèle de régression. Cette partie est présentée dans le chapitre 4 et l'annexe A.
 - Dans le chapitre 4, les méthodes d'estimations et les algorithmes utilisés dans le package sont décrits. Les performances des fonctions du package en termes de qualité prédictive de l'estimateur et d'estimation des indices de Sobol, sont validées par une étude de simulation.
 - Dans l'annexe A, la documentation complète du package, y compris des explications détaillées des fonctions du package ainsi que des exemples d'utilisation de chaque fonction du package, est fournie.

Les résumés des chapitres 3 et 4 sont présentés respectivement aux sections 1.2 et 1.3. Auparavant, plusieurs outils communs à ces deux chapitres sont brièvement décrits. Plus précisément:

- introduction sur l'analyse de sensibilité (voir la section 1.1.1),
- focus sur les méthodes basées sur la décomposition de la variance (voir la section 1.1.2),
- introduction sur la méta-modélisation (voir la section 1.1.3),
- construction d'un méta-modèle par projection sur des espaces à noyaux auto-reproduisants (RKHS) (voir la section 1.1.4),
- méthode d'estimation (voir la section 1.1.5).

1.1.1 Introduction sur l'analyse de sensibilité

Les méthodes d'analyse de sensibilité permettent d'étudier les relations entre les variables d'entrée et de sortie du modèle, et de mesurer l'effet de chaque variable ou

groupe de variables sur la sortie du modèle. Les objectifs principaux de l'analyse de sensibilité sont la calibration et la validation des modèles ainsi que l'aide à la prise de décision. Les méthodes et objectifs classiques de l'analyse de sensibilité sont décrits dans les ouvrages suivants: [Cacuci \(2003\)](#), [Fang et al. \(2005\)](#), [Dean and Lewis \(2006\)](#), [de Rocquigny et al. \(2008\)](#), [Saltelli \(2008\)](#), [Helton \(2008\)](#), [Saltelli et al. \(2009\)](#), [Faivre et al. \(2013\)](#), [Borgonovo and Plischke \(2016\)](#).

L'analyse de sensibilité repose sur le calcul et l'analyse des mesures qui évaluent l'effet des variables d'entrée sur la sortie du modèle. Par exemple, l'effet d'une variable d'entrée sur la sortie du modèle peut être évalué par la contribution de cette variable d'entrée sur la variance de la sortie du modèle. Les méthodes d'analyse de sensibilité peuvent être classées en deux groupes principaux:

L'analyse de sensibilité locale où il s'agit d'étudier l'impact local des variables d'entrée sur la variable de sortie. Elle consiste à calculer le gradient de la variable de sortie par rapport aux variables d'entrée autour d'une valeur choisie (la valeur moyenne des variables d'entrée par exemple). De nombreuses méthodes ont été développées pour calculer efficacement le gradient, notamment la modélisation Adjointe ([Cacuci \(2003\)](#), [Cacuci and Navon \(2005\)](#)) et la Différenciation Automatisée ([Griewank and Walther \(2008\)](#)). Les méthodes locales n'explorent pas pleinement l'espace des variables d'entrée, mais étudient l'impact de petites perturbations des variables d'entrée (généralement une variable à la fois) sur la variable de sortie.

L'analyse de sensibilité globale où il s'agit de calculer l'incertitude de la variable de sortie due aux variations des variables d'entrée. Contrairement à l'analyse de sensibilité locale, cette classe de méthodes prend en compte toute la gamme de variation des variables d'entrée. Les méthodes d'analyse de sensibilité globale sont nombreuses, voir par exemple [Saltelli et al. \(2009\)](#) pour un aperçu et [Iooss and Lemaître \(2015\)](#) pour une revue de ces méthodes. Généralement, les méthodes de l'analyse de sensibilité globale qui permettent de calculer les mesures de sensibilité quantitatives les plus utilisées peuvent être présentées en deux groupes:

- ✓ Les méthodes basées sur la régression sont appropriées lorsque le modèle est linéaire, c'est-à-dire si le coefficient de détermination R^2 est proche de un. Les mesures de sensibilité les plus utilisées dans ce cas sont: les coefficients de régression standardisés, les coefficients de corrélation de Pearson, et les coefficients de corrélation partielle. Dans le cas d'un modèle non-linéaire monotone, ces coefficients sont utilisés pour calculer des mesures de sensibilité, après avoir appliqué une transformation en rangs ([Saltelli et al. \(2009\)](#)). Lorsque le modèle est non-linéaire et non-monotone, ces méthodes ne produisent pas de mesures de sensibilité satisfaisantes ([Saltelli and Sobol \(1995\)](#)).
- ✓ Les méthodes basées sur la décomposition de la variance s'appliquent aux modèles non-linéaires et non-monotones. Il s'agit alors d'effectuer une décomposition de la variance de la variable de sortie. Plus précisément, la variance

de la variable de sortie est décomposée en parties attribuables à chacune des variables d'entrée et à leurs interactions. Les mesures de sensibilité sont exprimées comme le rapport de la variance due à chacun de ces groupes de variables (variables individuelles ou interactions de plusieurs variables) à la variance de la variable de sortie. La décomposition de la variance est pertinente si les variables d'entrée sont indépendantes les unes des autres (Saltelli and Tarantola (2002)). Ces méthodes sont largement utilisées car elles permettent d'explorer complètement l'espace des variables d'entrée, en tenant compte des effets d'interactions des variables d'entrée sur le modèle et de la non-linéarité du modèle.

1.1.2 Analyse de sensibilité globale: méthodes basées sur la décomposition de la variance

Considérons un modèle m dépendant de d variables d'entrée $X = (X_1, \dots, X_d)$ qui sont indépendantes et ont une loi connue $P_X = P_1 \otimes \dots \otimes P_d$ sur $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$, un sous-ensemble de \mathbb{R}^d . Le modèle m de \mathbb{R}^d vers \mathbb{R} est de carré-intégrable sur \mathcal{X} , c'est-à-dire $m \in L^2(\mathcal{X}, P_X)$.

Dans le cadre classique de l'analyse de sensibilité, où pour chaque valeur de X on peut calculer $m(X)$, on peut utiliser la méthode de Sobol (1993) brièvement décrite ci-dessous.

L'indépendance entre les coordonnées de X permet d'écrire le modèle m selon sa décomposition de Hoeffding (Sobol (1993), van der Vaart (1998)):

$$m(X) = m_0 + \sum_{a=1}^d m_a(X_a) + \sum_{a < a'} m_{a,a'}(X_a, X_{a'}) + \dots + m_{1,\dots,d}(X). \quad (1.2)$$

Les termes de cette décomposition sont définis en terme d'espérance conditionnelle:

$$\begin{aligned} m_0 &= E_X(m(X)), \\ m_a(X_a) &= E_X(m(X)|X_a) - m_0, \\ m_{a,a'}(X_a, X_{a'}) &= E_X(m(X)|X_a, X_{a'}) - m_a(X_a) - m_{a'}(X_{a'}) - m_0, \end{aligned}$$

et ainsi de suite pour les interactions d'ordre supérieur à deux.

Ces termes sont appelés terme constant, effets principaux, interactions d'ordre deux et d'ordre supérieur.

Soit \mathcal{P} l'ensemble de tous les sous-ensembles de $\{1, \dots, d\}$ de dimension 1 à d . Pour un ensemble A on note $|A|$ son cardinal. Pour tout $v \in \mathcal{P}$ et $X \in \mathcal{X}$, soit X_v le vecteur de composantes X_a , $a \in v$ et $m_v : \mathbb{R}^{|v|} \rightarrow \mathbb{R}$ la fonction associée à X_v dans l'équation (1.2). L'équation (1.2) peut alors être exprimée comme suit:

$$m(X) = m_0 + \sum_{v \in \mathcal{P}} m_v(X_v). \quad (1.3)$$

Cette décomposition est unique, tous les termes m_v , $v \in \mathcal{P}$, sont centrés et orthogonaux par rapport à $L^2(\mathcal{X}, P_X)$, c'est-à-dire,

$$\forall v \in \mathcal{P}, E_X(m_v(X_v)) = 0,$$

et

$$\forall v, v' \in \mathcal{P}, v \neq v', E_X(m_v(X_v)m_{v'}(X_{v'})) = 0.$$

Étant donné que les termes de la décomposition (1.3) sont centrés, de carré-intégrable et orthogonaux deux à deux par rapport à la distribution de X , la variance de $m(X)$ se décompose comme suit:

$$\text{var}(m(X)) = \sum_{v \in \mathcal{P}} \text{var}(m_v(X_v)). \quad (1.4)$$

Pour tout groupe de variables X_v , $v \in \mathcal{P}$, les indices de Sobol sont définis par:

$$S_v = \frac{\text{var}(m_v(X_v))}{\text{var}(m(X))}.$$

Pour chaque v , S_v exprime la fraction de la variance de $m(X)$ expliquée par X_v .

Pour tout $v \in \mathcal{P}$, quand $|v| = 1$, les S_v sont appelés indices du premier ordre ou indices des effets principaux. Quand $|v| = 2$, c'est-à-dire $v = \{a, a'\}$ et $a \neq a'$, ils sont appelés indices du second ordre ou indices d'interaction d'ordre deux (entre X_a et $X_{a'}$). Et ainsi de suite pour $|v| > 2$.

Le nombre total des indices de Sobol à calculer est égal à $|\mathcal{P}| = 2^d - 1$, qui augmente exponentiellement avec le nombre de variables d'entrée d . Lorsque d est grand, l'évaluation de tous les indices peut être très coûteuse voire même impossible. Pour cette raison, seuls les indices d'ordre inférieur ou égale à deux sont calculés en pratique. Cependant, les indices du premier et du second ordre ne peuvent pas toujours fournir une information sur la sensibilité du modèle. Afin de fournir une meilleure information sur la sensibilité du modèle, [Homma and Saltelli \(1996\)](#) ont proposé de calculer les indices du premier ordre et les indices d'ordre total définis comme suit:

Soit $\mathcal{P}_a \subset \mathcal{P}$ l'ensemble de tout les sous-ensembles de $\{1, \dots, d\}$ incluant a , alors

$$S_{T_a} = \sum_{v \in \mathcal{P}_a} S_v.$$

Pour tout $a \in \{1, \dots, d\}$, S_{T_a} indique l'effet total de la variable X_a . Il exprime la fraction de la variance expliquée par la variable X_a seule et toute interaction de X_a avec les autres variables.

Les indices d'ordre total permettent de classer les variables d'entrée selon la quantité de leur effet sur la variable de sortie. Néanmoins, ils ne fournissent pas d'informations complètes sur la sensibilité du modèle comme le font tous les indices de Sobol.

Le calcul des indices de Sobol est généralement effectué par les méthodes de Monte Carlo (voir par exemple: [Sobol \(1993\)](#) pour les effets principaux et interactions, et [Saltelli \(2002\)](#) pour les effets principaux et indices d'ordre total). Ces méthodes sont très coûteuses, car elles peuvent nécessiter le calcul du modèle plusieurs milliers de fois pour obtenir des estimations précises des indices de Sobol. Ainsi dans

le cas où d est grand, m est complexe et où le calcul des variances est numériquement compliqué voire impossible, comme dans le cas où le modèle m est inconnu, les méthodes décrites ci-dessus ne sont pas pertinentes.

Une autre méthode consiste à approcher m par un modèle simplifié, appelé méta-modèle, qui est beaucoup plus rapide à évaluer, et à effectuer l'analyse de sensibilité sur celui-ci. Non seulement un méta-modèle permet de calculer à moindre coût des indices de Sobol approchés, mais il fournit de l'information sur la nature des effets des variables d'entrées ou de leurs interactions sur la variable de sortie.

1.1.3 Méta-modélisation

La méta-modélisation consiste à construire une fonction qui est calculable, facile à interpréter et qui a de bonnes qualités de prédiction. Soit $\{m(X_i)\}_{i=1}^n$ les résultats de n évaluations du modèle m basées sur un plan d'expérience $\{X_i\}_{i=1}^n$. Dans ce contexte, un méta-modèle est une approximation du modèle m construite à partir du plan d'expérience $\{X_i\}_{i=1}^n$ et des sorties $\{m(X_i)\}_{i=1}^n$. Il existe différentes approches de méta-modélisation, voir [Sacks et al. \(1989\)](#), [Friedman \(1991\)](#), [Breiman \(2001\)](#), [Friedman \(2001\)](#), [Kennedy and O'Hagan \(2001\)](#), [Oakley and O'Hagan \(2004\)](#), [Storlie and Helton \(2008\)](#), [Storlie et al. \(2009\)](#), [Storlie et al. \(2011\)](#) pour différents exemples, et [Touzani \(2011\)](#) pour un aperçu.

Dans le cadre de l'analyse de sensibilité globale, on considère un méta-modèle dont la décomposition additive est candidate pour approcher la décomposition de Hoeffding de m . Ce méta-modèle permettra ainsi d'effectuer l'analyse de sensibilité globale de m , en calculant des indices de Sobol, éventuellement d'ordre élevé. Par une fonction qui a la décomposition additive, on entend une fonction f de $\mathcal{X} \subset \mathbb{R}^d$ vers \mathbb{R} qui est définie comme suit:

$$f = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \quad E_X(f_v(X_v)) = 0, \quad E_X(f_v(X_v)f_{v'}(X_{v'})) = 0, \quad \forall v, v' \in \mathcal{P}, \quad v \neq v',$$

où f_0 est une constante, et les fonctions f_v sont supposées appartenir aux espaces fonctionnels.

Parmi les approches de méta-modélisation proposées dans la littérature, la décomposition basée sur les polynômes de Chaos ([Wiener \(1938\)](#), [Schoutens \(2000\)](#)) permet d'approcher la décomposition de Hoeffding de m ([Sudret \(2008\)](#)).

Le principe de la décomposition selon les polynômes de Chaos est de projeter m sur une base de polynômes orthonormés de la façon suivante ([Soize and Ghanem \(2004\)](#)):

$$m(X) = \sum_{j=0}^{\infty} h_j \phi_j(X), \tag{1.5}$$

où $\{h_j\}_{j=0}^{\infty}$ sont les coefficients, et $\{\phi_j\}_{j=0}^{\infty}$ sont des polynômes orthonormés multivariés associés à X qui sont déterminés par la distribution des coordonnées de X . En pratique, la série définie en (1.5) doit être tronquée conduisant à approcher m

par:

$$m(X) \approx \sum_{j=0}^{v_{max}} h_j \phi_j(X), \quad (1.6)$$

où v_{max} doit être déterminé par une méthode numérique. Dans cette approche, les indices de Sobol sont explicitement donnés à partir des carrés des coefficients associés.

[Blatman and Sudret \(2011\)](#) ont proposé une méthode pour tronquer la série (1.5) et un algorithme basé sur la méthode *least-angle regression* pour sélectionner les termes pertinents dans le développement.

Dans cette approche, la famille des polynômes orthonormés $\{\phi_j\}_{j=0}^{\infty}$ est déterminée de manière unique par la distribution des coordonnées de X . Cependant, cette famille ne constitue pas nécessairement la meilleure base fonctionnelle pour bien approcher m .

Une autre approche pour construire des méta-modèles est la modélisation par le processus gaussien (GP) ([Welch et al. \(1992\)](#), [Oakley and O’Hagan \(2004\)](#), [Kleijnen \(2007, 2009\)](#), [Marrel et al. \(2009\)](#), [Durrande et al. \(2012\)](#), [Le Gratiet et al. \(2014\)](#)). Le principe est de modéliser la distribution a priori de $m(X)$ par un modèle de GP, noté $\mathcal{Z}(X)$, de moyenne $m_{\mathcal{Z}}(X)$ et de noyau de covariance $k_{\mathcal{Z}}(X, X')$. Pour effectuer l’analyse de sensibilité, on peut remplacer le vrai modèle $m(X)$ par l’espérance de la loi à posteriori de $\mathcal{Z}(X)$, et en déduire les indices de Sobol. La plupart du temps, avec GP, les indices de Sobol sont estimés à l’aide de méthodes de Monte Carlo.

Une revue de la méta-modélisation basée sur les polynômes de Chaos et le GP est présentée dans l’ouvrage de [Le Gratiet et al. \(2017\)](#).

[Durrande et al. \(2013\)](#) ont considéré une classe de méthodes d’approximation fonctionnelle similaire au GP et ont obtenu un méta-modèle qui satisfait les propriétés de la décomposition de Hoeffding. Ils ont proposé d’approcher m par des fonctions appartenant à un RKHS \mathcal{H} qui est construit comme une somme directe d’espaces RKHS, de sorte que la projection de m sur \mathcal{H} est une approximation de la décomposition de Hoeffding de m .

Dans le modèle de régression, lorsque les valeurs de $\{m(X_i)\}_{i=1}^n$ ne peuvent pas être calculées, on peut utiliser les méthodes de projection sur une base fonctionnelle pour estimer un méta-modèle pour m . Ce méta-modèle sera estimé en utilisant des approches d’estimation non-paramétriques à partir des observations $\{(X_i, Y_i)\}_{i=1}^n$, et on déduira de cet estimateur des estimateurs des indices de Sobol de m . [Huet and Taupin \(2017\)](#) ont considéré les mêmes espaces d’approximation fonctionnels que [Durrande et al. \(2013\)](#), et ont proposé un estimateur d’un méta-modèle qui approche la décomposition de Hoeffding de m . Elles ont déduit de ce méta-modèle estimé, des estimateurs pour les indices de Sobol de m . Cette approche est présentée plus en détail par la suite.

1.1.4 Méta-modèles basés sur des espaces à noyaux auto-reproduisants (RKHS)

Cette section commencera tout d'abord par une brève introduction des espaces RKHS. La méthode de [Durrande et al. \(2013\)](#) pour construire le RKHS et la définition de méta-modèle f^* qui approche la décomposition de Hoeffding de m , sont décrites respectivement dans les sections 1.1.4.2 et 1.1.4.3.

1.1.4.1 Introduction aux espaces RKHS

Soit \mathcal{H} un espace de Hilbert de fonctions définies sur un ensemble \mathcal{X} . L'espace \mathcal{H} est un RKHS si pour tout $X \in \mathcal{X}$ les fonctionnelles,

$$\begin{aligned} L_X : \mathcal{H} &\rightarrow \mathbb{R} \\ f &\mapsto f(X), \end{aligned}$$

sont continues.

Le théorème de représentation de Riesz assure l'existence d'un élément unique $k_X(\cdot)$ dans \mathcal{H} vérifiant la propriété suivante:

$$\forall X \in \mathcal{X}, \forall f \in \mathcal{H}, f(X) = L_X(f) = \langle f, k_X \rangle_{\mathcal{H}}.$$

où $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ dénote le produit scalaire dans \mathcal{H} .

Il en découle que pour tout X, X' dans \mathcal{X} , et $k_X(\cdot), k_{X'}(\cdot)$ dans \mathcal{H} , on a,

$$k_X(X') = L_{X'}(k_X) = \langle k_X, k_{X'} \rangle_{\mathcal{H}}. \quad (1.7)$$

Cela permet de définir le noyau auto-reproduisant de \mathcal{H} comme suit:

$$\begin{aligned} k : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ (X, X') &\mapsto k_X(X'). \end{aligned}$$

Le noyau auto-reproduisant $k(X, X')$ satisfait les propriétés suivantes:

- Il est symétrique. En effet, par définition de $k(\cdot, \cdot)$ et grâce à la propriété (1.7), on a:

$$k(X, X') = k_X(X') = \langle k_X, k_{X'} \rangle_{\mathcal{H}} = k_{X'}(X) = k(X', X).$$

- Pour tout $n \in \mathbb{N}$, $\{X_i\}_{i=1}^n \in \mathcal{X}$ et $\{c_i\}_{i=1}^n \in \mathbb{R}$, on a:

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(X_i, X_j) &= \sum_{i=1}^n \sum_{j=1}^n \langle c_i k(X_i, \cdot), c_j k(X_j, \cdot) \rangle_{\mathcal{H}}, \\ &= \left\| \sum_{i=1}^n c_i k(X_i, \cdot) \right\|_{\mathcal{H}}^2 \geq 0. \end{aligned}$$

Ainsi, $k(X, X')$ est défini positif.

Davantage d'informations sur les espaces RKHS sont indiquées dans des ouvrages standards comme [Aronszajn \(1950\)](#), [Saitoh \(1988\)](#) et [Berlinet and Thomas-Agnan \(2003\)](#).

1.1.4.2 Construction du RKHS et décomposition de Hoeffding

L'idée est de construire un RKHS incluant les fonctions qui ont la décomposition additive et qui sont candidates pour approcher la décomposition de Hoeffding de m . Pour cela, on utilise la méthode de [Durrande et al. \(2013\)](#) décrite ci-dessous.

Soit $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ un sous-ensemble de \mathbb{R}^d . Pour chaque $a \in \{1, \dots, d\}$, on choisit un RKHS \mathcal{H}_a et son noyau associé k_a défini sur l'ensemble $\mathcal{X}_a \subset \mathbb{R}$, tels que les deux propriétés suivantes soient satisfaites:

- (i) $k_a : \mathcal{X}_a \times \mathcal{X}_a \rightarrow \mathbb{R}$ est $P_a \otimes P_a$ mesurable,
- (ii) $E_{X_a} \sqrt{k_a(X_a, X_a)} < \infty$.

La propriété (ii) dépend du noyau k_a , $a = 1, \dots, d$ et de la loi de X_a , $a = 1, \dots, d$. Elle est relativement peu restrictive car elle est satisfaite, par exemple, pour tous les noyaux bornés.

Le RKHS \mathcal{H}_a peut être décomposé en une somme de deux sous-RKHS orthogonaux,

$$\mathcal{H}_a = \mathcal{H}_{0a} \overset{\perp}{\oplus} \mathcal{H}_{1a},$$

où \mathcal{H}_{0a} est le RKHS des fonctions centrées,

$$\mathcal{H}_{0a} = \left\{ f_a \in \mathcal{H}_a : E_{X_a}(f_a(X_a)) = 0 \right\},$$

et \mathcal{H}_{1a} est le RKHS des fonctions constantes,

$$\mathcal{H}_{1a} = \left\{ f_a \in \mathcal{H}_a : f_a(X_a) = C \right\}.$$

Le noyau k_{0a} associé au RKHS \mathcal{H}_{0a} est défini par:

$$k_{0a}(X_a, X'_a) = k_a(X_a, X'_a) - \frac{E_{U \sim P_a}(k_a(X_a, U))E_{U \sim P_a}(k_a(X'_a, U))}{E_{(U, V) \sim P_a \otimes P_a} k_a(U, V)}. \quad (1.8)$$

Soit $k_v(X_v, X'_v) = \prod_{a \in v} k_{0a}(X_a, X'_a)$, le noyau ANOVA $k(\cdot, \cdot)$ est défini comme suit:

$$k(X, X') = \prod_{a=1}^d \left(1 + k_{0a}(X_a, X'_a) \right) = 1 + \sum_{v \in \mathcal{P}} k_v(X_v, X'_v).$$

Pour \mathcal{H}_v étant le RKHS associé au noyau k_v , le RKHS associé au noyau ANOVA est défini par,

$$\mathcal{H} = \prod_{a=1}^d \left(\mathbb{1} \overset{\perp}{\oplus} \mathcal{H}_{0a} \right) = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v,$$

où \perp correspond à une orthogonalité pour le produit scalaire sur L^2 .

D'après cette construction, toute fonction $f \in \mathcal{H}$ satisfait la décomposition suivante:

$$f(X) = \langle f, k(X, \cdot) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v),$$

qui est la décomposition de Hoeffding de f .

Les propriétés de régularité du RKHS \mathcal{H} construit comme décrit ci-dessus, dépendent de l'ensemble des noyaux $(k_a, a = 1, \dots, d)$. Cette méthode permet de choisir différents espaces d'approximation indépendamment de la distribution des coordonnées de X , en choisissant différents ensembles de noyaux. Alors que, comme indiqué précédemment, dans l'approche de méta-modélisation basée sur le développement en polynômes de Chaos, la famille des polynômes orthonormés $\{\phi_j\}_{j=0}^\infty$ est déterminée de manière unique par la distribution des coordonnées de X . Ici, la distribution des coordonnées de X n'intervient que pour l'orthogonalisation des espaces \mathcal{H}_v , $v \in \mathcal{P}$ mais pas dans le choix des RKHS, pourvu que les propriétés (i) et (ii) soient satisfaites. C'est l'un des principaux avantages de cette méthode par rapport à l'approche basée sur le développement en polynômes de Chaos où la régularité de l'approximation n'est gérée que par le choix de v_{max} (voir l'équation (1.6)) et non par celui de la base fonctionnelle (Blatman and Sudret (2011)).

1.1.4.3 Approximation de la décomposition de Hoeffding de m

Soit f^* la projection orthogonale de m sur \mathcal{H} définie par:

$$f^* = \arg \min_{f \in \mathcal{H}} \|m - f\|_2^2 = \arg \min_{f \in \mathcal{H}} E_X (m(X) - f(X))^2.$$

La fonction $f^* \in \mathcal{H}$, $f^* = f_0^* + \sum_{v \in \mathcal{P}} f_v^*$ est l'approximation de m sur le RKHS \mathcal{H} , et sa décomposition de Hoeffding est une approximation de la décomposition de Hoeffding de m . Par conséquent, pour chaque $v \in \mathcal{P}$, la fonction f_v^* approche la fonction m_v dans l'équation (1.3).

Le nombre de fonctions f_v^* est lié au cardinal de \mathcal{P} , égal à $2^d - 1$, qui peut devenir très grand dès que d est grand. Ainsi, l'idée est de calculer un estimateur sparse de f^* comme estimateur de m en utilisant des méthodes d'estimation non-paramétriques.

1.1.5 Méthode d'estimation

Considérons le modèle de régression défini dans l'équation (1.1),

$$Y = m(X) + \sigma\varepsilon, \sigma > 0.$$

La fonction inconnue m est approchée par le méta-modèle f^* qui est ensuite estimé par un estimateur sparse \hat{f} . Cet estimateur \hat{f} , basé sur n observations $\{(Y_i, X_i)\}_{i=1}^n$, minimise un critère pénalisé. La fonction de pénalité prend en compte à la fois la nature non-paramétrique du problème et le nombre éventuellement important de fonctions qui doivent être estimées.

Avant de décrire la méthode pour calculer \hat{f} , on rappelle quelques méthodes liées à l'estimation dans un modèle de régression non-paramétrique additif.

Certains auteurs approchent m par une fonction qui a une décomposition additive univariée et sparse de la forme,

$$f(X) = f_0 + \sum_{a \in S} f_a(X_a) \text{ avec } |S| < d, \quad (1.9)$$

où f_0 est une constante et où pour tout $a \in S$, les f_a sont des fonctions supposées régulières. Les fonctions f_a sont estimées à partir des observations à l'aide d'un critère pénalisé.

Ravikumar et al. (2009) ont considéré un espace de Hilbert \mathcal{H} de fonctions qui ont une forme additive univariée. Leur espace d'approximation fonctionnelle \mathcal{H} est construit comme une somme directe des espaces de Hilbert, c'est-à-dire

$$\mathcal{H} = \bigoplus_{a=1}^d \mathcal{H}_a,$$

où pour tout $a \in \{1, \dots, d\}$, \mathcal{H}_a est le sous-espace de Hilbert de $L^2(\mathcal{X}_a, P_a)$ des fonctions univariées f_a qui sont centrées et P_a mesurables. Afin de favoriser sparsité et régularité, ils ont proposé la méthode SpAM (*Sparse Additive Models*). Leur méthode est basée sur la minimisation du critère des moindres carrés pénalisé par une fonction de pénalité définie comme suit,

$$\lambda \sum_{a=1}^d \sqrt{\int (f_a(X_a))^2 dX_a}, \lambda \in \mathbb{R}^+.$$

Meier et al. (2009) ont proposé un estimateur qui est dans l'espace des splines cubiques naturel. Leur méthode est basée sur la minimisation du critère des moindres carrés pénalisé par une fonction de pénalité de la forme,

$$\sum_{a=1}^d \sqrt{\lambda_1 \|f_a\|_n^2 + \lambda_2 \int (f_a''(X_a))^2 dX_a}, \lambda_1, \lambda_2 \in \mathbb{R}^+,$$

où $\|f_a\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_a^2(X_{ai})$.

Leur fonction de pénalité est composée de deux parties: la première partie favorise la sparsité et la deuxième partie favorise la régularité.

Raskutti et al. (2012) ont considéré plusieurs espaces d'approximation fonctionnelle, y compris les polynômes, les splines et les classes de Sobolev. Leur méthode est basée sur la minimisation du critère des moindres carrés pénalisé par une fonction de pénalité de la forme,

$$\gamma \|f\|_{n,1} + \mu \|f\|_{\mathcal{H},1}, \gamma, \mu \in \mathbb{R}^+, \quad (1.10)$$

où $\|f\|_{n,1} = \sum_{a=1}^d \|f_a\|_n$ et $\|f\|_{\mathcal{H},1} = \sum_{a=1}^d \|f_a\|_{\mathcal{H}_a}$.

Dans leur fonction de pénalité, la première partie favorise la sparsité et la deuxième partie favorise la régularité.

Effectuer l'analyse de sensibilité globale sur un modèle additif univarié conduit à n'obtenir que les indices de Sobol de premier ordre, ce qui ne fournit peut-être pas une bonne information sur la sensibilité du modèle. Les interactions entre les variables qui peuvent affecter la relation entre Y et X sont complètement ignorées dans ce contexte.

Afin d'inclure les effets d'interaction, on peut approcher m par une fonction qui a une décomposition additive multivariée et sparse de la forme,

$$f(X) = f_0 + \sum_{v \in S} f_v(X_v) \text{ avec } |S| < |\mathcal{P}|,$$

qui est une généralisation de la décomposition additive univariée définie dans l'équation (1.9).

Dans le cadre du lissage par splines de type ANOVA (Wahba (1990), Friedman (1991), Wahba et al. (1995)), Lin and Zhang (2006) ont proposé la méthode COSSO (*Component Selection and Smoothing Operator*). Leur méthode est basée sur la minimisation du critère des moindres carrés pénalisé par une fonction de pénalité qui est la combinaison de la norme l_1 avec la norme de Hilbert. L'implémentation de COSSO s'effectue sur les espaces de Sobolev de second ordre.

Kandasamy and Yu (2016) ont proposé la méthode SLASA (*Shrunk Additive Least Squares Approximation*) qui est basée sur la minimisation du critère des moindres carrés pénalisé par la somme des carrés des normes RKHS. Leur estimateur est une fonction additive multivariée d'ordre v_{max} contenant $\binom{d}{v_{max}}$ termes dans son développement. La valeur de v_{max} est déterminée en utilisant une procédure de validation croisée.

Huet and Taupin (2017) ont considéré un estimateur d'un méta-modèle qui approche la décomposition de Hoeffding de m définie dans l'équation (1.3). Leur estimateur est la solution de minimisation du critère des moindres carrés pénalisé, où la fonction de pénalité est définie dans l'équation (1.10) et est adaptée au cadre multivarié,

$$\gamma \|f\|_n + \mu \|f\|_{\mathcal{H}} \text{ avec } \|f\|_{\mathcal{H}} = \sum_{v \in \mathcal{P}} \|f_v\|_{\mathcal{H}_v}, \text{ et } \|f\|_n = \sum_{v \in \mathcal{P}} \|f_v\|_n.$$

Leur méthode, appelée *ridge group sparse*, estime les groupes v qui sont pertinents pour prédire le méta-modèle f^* et la relation entre f_v^* et X_v pour chaque groupe $v \in \mathcal{P}$. L'estimateur obtenu, appelé l'estimateur *RKHS ridge group sparse*, est ensuite utilisé pour estimer les indices de Sobol de m . Cette méthode permet d'estimer les indices de Sobol pour tous les groupes dans le support de l'estimateur *RKHS ridge group sparse*, y compris les interactions d'ordre élevé, un point connu pour être difficile à mettre en pratique.

Dans ce travail, la méthode proposée par Huet and Taupin (2017) est utilisée afin de calculer un estimateur sparse du méta-modèle f^* qui permet également de calculer les estimateurs des indices de Sobol de m . Décrivons plus en détail cette méthode et la méthode pour estimer les indices de Sobol de m respectivement dans la section suivante et la section 1.1.5.2.

1.1.5.1 Procédure *ridge group sparse* et estimateur associé

Pour tout $v \in \mathcal{P}$, soit X_v la matrice des variables correspondant au v -ième groupe,

$$X_v = (X_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}.$$

Pour tout $f \in \mathcal{H}$ tel que $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, et pour des paramètres de régularisations $\gamma_v, \mu_v, v \in \mathcal{P}$, le critère *ridge group sparse* est défini comme suit:

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \left(Y_i - f_0 - \sum_{v \in \mathcal{P}} f_v(X_{vi}) \right)^2 + \sum_{v \in \mathcal{P}} \gamma_v \|f_v\|_n + \sum_{v \in \mathcal{P}} \mu_v \|f_v\|_{\mathcal{H}_v}, \quad (1.11)$$

où $\|f_v\|_n$ est la norme empirique L^2 de f_v définie en fonction de l'échantillon $\{X_{vi}\}_{i=1}^n$ comme suit:

$$\|f_v\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_v^2(X_{vi}).$$

La fonction de pénalité dans le critère $\mathcal{L}(f)$ est la somme de la norme empirique et de la norme de Hilbert, ce qui permet de sélectionner peu de termes dans la décomposition additive de f sur les ensembles $v \in \mathcal{P}$. De plus, la norme de Hilbert favorise la régularité du $f_v, v \in \mathcal{P}$ estimé.

Définissons l'ensemble des fonctions,

$$\mathcal{F} = \left\{ f : f = f_0 + \sum_{v \in \mathcal{P}} f_v, \text{ with } f_v \in \mathcal{H}_v, \text{ and } \|f_v\|_{\mathcal{H}_v} \leq r_v, r_v > 0 \right\}. \quad (1.12)$$

L'estimateur *RKHS ridge group sparse* de m est défini par:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f). \quad (1.13)$$

D'après le *representer théorème* (Kimeldorf and Wahba (1970)), le problème de minimisation fonctionnelle non-paramétrique (1.13) est équivalent à un problème de minimisation paramétrique. En effet, la solution du problème de minimisation (1.13) appartenant au RKHS \mathcal{H} est écrite comme $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, où pour une matrice $\theta = (\theta_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}$ on a pour tout $v \in \mathcal{P}$,

$$f_v(\cdot) = \sum_{i=1}^n \theta_{vi} k_v(X_{vi}, \cdot).$$

Soit $\|\cdot\|$ la norme euclidienne dans \mathbb{R}^n , et pour chaque $v \in \mathcal{P}$, soit K_v la $n \times n$ matrice de Gram associée au noyau $k_v(\cdot, \cdot)$, c'est-à-dire

$$(K_v)_{i,i'} = k_v(X_{vi}, X_{vi'}).$$

Soit aussi $K_v^{1/2}$ la matrice qui satisfait $t(K_v^{1/2})K_v^{1/2} = K_v$, et soit \hat{f}_0 et $\hat{\theta}$ les solutions de minimisation du critère suivant:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \sum_{v \in \mathcal{P}} \gamma_v \|K_v \theta_v\| + n \sum_{v \in \mathcal{P}} \mu_v \|K_v^{1/2} \theta_v\|. \quad (1.14)$$

Alors l'estimateur \hat{f} défini dans l'équation (1.13) satisfait,

$$\hat{f}(X) = \hat{f}_0 + \sum_{v \in \mathcal{P}} \hat{f}_v(X_v) \text{ avec } \hat{f}_v(X_v) = \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v). \quad (1.15)$$

Comme le critère $C(f_0, \theta)$ est convexe et séparable, on peut calculer $\hat{\theta}$ en utilisant un algorithme de *block coordinate descent* (Boyd et al. (2011), Bubeck (2015)).

Remarque 1.1.1 *La contrainte $\|f_v\|_{\mathcal{H}_v} \leq r_v$ n'est pas prise en compte dans le problème de minimisation paramétrique. Cette contrainte est cruciale pour les propriétés théoriques, mais la valeur de r_v est inconnue et n'est pas utile en pratique.*

1.1.5.2 Estimation des indices de Sobol de m

La variance de la fonction m est estimée par la variance de l'estimateur \hat{f} . Comme l'estimateur \hat{f} appartient au RKHS \mathcal{H} , il admet la décomposition de Hoeffding et,

$$\text{var}(\hat{f}(X)) = \sum_{v \in \mathcal{P}} \text{var}(\hat{f}_v(X_v)),$$

où pour tout $v \in \mathcal{P}$,

$$\text{var}(\hat{f}_v(X_v)) = E_X(\hat{f}_v^2(X_v)) = \|\hat{f}_v\|_2^2.$$

Afin de réduire le temps de calcul, on peut estimer les variances de $\hat{f}_v(X_v)$, $v \in \mathcal{P}$ par leurs variances empiriques.

Soit \hat{f}_v la moyenne empirique de $\{\hat{f}_v(X_{vi})\}_{i=1}^n$,

$$\hat{f}_v = \frac{1}{n} \sum_{i=1}^n \hat{f}(X_{vi}),$$

alors

$$\widehat{\text{var}}(\hat{f}_v(X_v)) = \frac{1}{n-1} \sum_{i=1}^n (\hat{f}_v(X_{vi}) - \hat{f}_v)^2.$$

Pour les groupes v qui appartiennent au support de \hat{f} , les estimateurs des indices de Sobol de m sont définis par,

$$\hat{S}_v = \frac{\widehat{\text{var}}(\hat{f}_v(X_v))}{\sum_{v \in \mathcal{P}} \widehat{\text{var}}(\hat{f}_v(X_v))},$$

et pour les groupes v qui n'appartiennent pas au support de \hat{f} , on a $\hat{S}_v = 0$.

1.2 Résumé du chapitre 3

1.2.1 Objectifs et résultats

Pour un estimateur \hat{f} d'un modèle m soit $R(m, \hat{f})$ son risque. Le risque $R(m, \hat{f})$ est une mesure qui caractérise la précision de l'estimateur \hat{f} et qui peut être exprimé en fonction du biais et de la variance de \hat{f} :

$$R(m, \hat{f}) = (\text{biais}(\hat{f}))^2 + \text{var}(\hat{f}).$$

Ainsi, la qualité de l'estimateur \hat{f} peut être mesurée par son risque. On considère le risque empirique L^2 de l'estimateur \hat{f} , c'est-à-dire lorsque $R(m, \hat{f}) = \|m - \hat{f}\|_n^2$, et le risque quadratique de l'estimateur \hat{f} , c'est-à-dire lorsque $R(m, \hat{f}) = \|m - \hat{f}\|_2^2$.

On s'intéresse à des propriétés non-asymptotiques de l'estimateur \hat{f} , au sens où l'on ne suppose pas que le nombre d'observations n tend vers l'infini. Nos résultats sont donc valables pour tout n avec une grande probabilité. On établit en particulier, des majorations du risque $R(m, \hat{f})$ de la forme,

$$R(m, \hat{f}) \leq C \inf_{f \in \mathcal{F}} \{R(m, f) + r_n(f)\}, \quad (1.16)$$

où C est une constante, et \mathcal{F} est l'espace d'approximation.

Soit f' la fonction dans \mathcal{F} pour laquelle l'infimum du membre de droite de l'inégalité (1.16) est réalisé. Le terme $R(m, f')$ est le terme de biais qui dépend du choix de l'espace d'approximation. Le terme $r_n(f)$ est le terme de variance qui doit décroître avec n . Il contrôle la vitesse de convergence, c'est-à-dire la vitesse à laquelle le risque de l'estimateur va s'approcher du meilleur possible. Le terme de variance dépend de la régularité des noyaux k_v , $v \in \mathcal{P}$, du nombre de termes intervenant dans la décomposition de la fonction f sur l'espace d'approximation, du nombre de variables d'entrée d , et du nombre d'observations n .

Dans le modèle de régression gaussienne, c'est-à-dire lorsque ε dans l'équation (1.1) est une variable gaussienne centrée, [Huet and Taupin \(2017\)](#) ont établi les majorations du risque empirique L^2 et du risque quadratique de l'estimateur *RKHS ridge group sparse* \hat{f} .

Dans ce chapitre, on considère le modèle de régression avec l'erreur ε non-gaussienne et non-bornée. Dans ce contexte, l'objectif est d'établir des majorations du risque de l'estimateur *RKHS ridge group sparse* \hat{f} de la forme (1.16) avec la même vitesse de convergence que dans le modèle de régression gaussienne. Les majorations du risque empirique L^2 et du risque quadratique de l'estimateur \hat{f} sont présentées respectivement dans le [résultat 1](#) et le [résultat 2](#).

1.2.2 Présentation du modèle

Considérons le modèle de régression défini dans l'équation (1.1),

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0.$$

Les variables d'entrée $X = (X_1, \dots, X_d)$ sont indépendantes et ont une loi connue $P_X = \bigotimes_{a=1}^d P_a$ sur $\mathcal{X} = \prod_{a=1}^d \mathcal{X}_a$, un sous-ensemble compact de \mathbb{R}^d . La fonction $m : \mathbb{R}^d \rightarrow \mathbb{R}$ est inconnue, peut-être complexe, et elle est supposée être de carré-intégrable sur \mathcal{X} .

Soit \mathcal{D} l'ensemble des densités,

$$\mathcal{D} = \left\{ \pi_\alpha : \pi_\alpha(x) = a_\alpha \exp(-|x|^\alpha), \text{ avec } (a_\alpha)^{-1} = \int_{\mathbb{R}} \exp(-|x|^\alpha) dx, \alpha > 2 \right\}. \quad (1.17)$$

Dans ce chapitre, on suppose que l'erreur ε est égale à Z/σ_α , où Z est une variable aléatoire de densité $\pi_\alpha \in \mathcal{D}$ et $\sigma_\alpha^2 = \text{var}(Z)$.

1.2.3 Les principaux résultats

Commençons par quelques notations, propriétés et hypothèses qui sont nécessaires pour annoncer le [résultat 1](#) et le [résultat 2](#).

Notations

✓ Pour une fonction $f \in \mathcal{H}$, soit S_f son support,

$$S_f = \{v \in \mathcal{P} : f_v \neq 0\}.$$

✓ À chaque noyau k_v , $v \in \mathcal{P}$ on associe l'opérateur intégral T_{k_v} de $L^2(\mathcal{X}_v, P_v)$ vers $L^2(\mathcal{X}_v, P_v)$ défini par:

$$\forall f \in L^2(\mathcal{X}_v, P_v), T_{k_v}(f) = \int_{\mathcal{X}_v} k_v(., t) f(t) dP_v(t).$$

Pour chaque $v \in \mathcal{P}$, soit $\omega_{v,1} \geq \omega_{v,2} \geq \dots \geq 0$ les valeurs propres de l'opérateur intégral T_{k_v} . Définissons la fonction $Q_{n,v}(t)$ pour un t positif comme suit:

$$Q_{n,v}(t) = \sqrt{\frac{5}{n} \sum_{\ell \geq 1} \min(t^2, \omega_{v,\ell})}.$$

✓ Pour une constante $\Delta > 0$, soit $\nu_{n,v}$ défini par:

$$\nu_{n,v} = \inf_t \left\{ Q_{n,v}(t) \leq \Delta t^2 \right\}. \quad (1.18)$$

Pour chaque $v \in \mathcal{P}$, $\nu_{n,v}$ est la vitesse minimax d'estimation par rapport à la norme $L^2(\mathcal{X}, P_X)$ dans le RKHS \mathcal{H}_v ([Mendelson \(2002\)](#)).

Remark 1.2.1 *la vitesse d'estimation $\nu_{n,v}$, $v \in \mathcal{P}$, est liée à la régularité du RKHS via le taux de décroissant des valeurs propres $\{\omega_{v,\ell}\}_{\ell=1}^\infty$. Lorsque le RKHS est très régulière, c'est-à-dire lorsque les valeurs propres $\{\omega_{v,\ell}\}_{\ell=1}^\infty$ tendent rapidement vers 0, la vitesse $\nu_{n,v}$, $v \in \mathcal{P}$ sera proche de la vitesse paramétrique (voir section 3.3.1 du chapitre 3).*

Propriétés

La construction du RKHS décrite à la section 1.1.4.2 assure que les propriétés suivantes sont satisfaites:

P1 Pour tout $v \in \mathcal{P}$, les fonctions $f_v \in \mathcal{H}_v$ sont centrées et de carré-intégrables,

$$E_X(f_v(X_v)) = 0 \text{ et } E_X(f_v^2(X_v)) < \infty.$$

P2 Pour tout $v, v' \in \mathcal{P}$, $v \neq v'$, les fonctions $f_v \in \mathcal{H}_v$ et $f_{v'} \in \mathcal{H}_{v'}$ sont orthogonales par rapport à $L^2(\mathcal{X}, P_X)$,

$$E_X(f_v(X_v) f_{v'}(X_{v'})) = 0.$$

Hypothèses

H1 Pour tout $v \in \mathcal{P}$, les fonctions $f_v \in \mathcal{H}_v$ sont uniformément bornées,

$$\exists R > 0 \text{ tel que } \|f_v\|_\infty = \sup_{X_v} |f_v(X_v)| \leq R.$$

Cette hypothèse est satisfaite lorsque le noyau k_v est borné sur l'ensemble compact \mathcal{X} . En effet,

$$\|f_v\|_\infty \leq \sup_{X \in \mathcal{X}} \sqrt{k_v(X_v, X_v)} \|f_v\|_{\mathcal{H}_v}.$$

Pour chaque $v \in \mathcal{P}$, soit $\lambda_{n,v}$ défini de la façon suivante:

$$\lambda_{n,v} = \max\left(\nu_{n,v}, \sqrt{\frac{d}{n}}\right). \quad (1.19)$$

Les paramètres de régularisation μ_v et γ_v intervenant dans le critère (1.11) sont choisis comme suit:

H2 Pour une constante $C_1 > 10 + 4\Delta$,

$$\forall v \in \mathcal{P}, \mu_v = C_1 \lambda_{n,v}^2, \gamma_v = C_1 \lambda_{n,v}.$$

H3 Il existe des constantes positives C_2, C_3 , et $0 < \beta < 1/\alpha$ telles que les conditions suivantes sont satisfaites:

$$\forall v \in \mathcal{P}, n \lambda_{n,v}^2 \geq -C_2 \log \lambda_{n,v}, \quad (1.20)$$

et

$$\forall f \in \mathcal{F}, \sum_{v \in S_f} \lambda_{n,v}^2 \leq C_3 n^{2\beta-1}. \quad (1.21)$$

Le chapitre 3 présente les deux résultats suivants.

Résultat 1: la majoration du risque empirique L^2 de l'estimateur *RKHS ridge group sparse*

Considérons le modèle de régression décrit dans la section 1.2.2 avec $\sigma = 1$. Soit $\{(Y_i, X_i)\}_{i=1}^n$ un échantillon de taille n de la même loi que (Y, X) , et soient $\{\varepsilon_i\}_{i=1}^n$ les erreurs aléatoires qui sont indépendantes et identiquement distribuées (i.i.d.) comme ε . Soit aussi l'estimateur *RKHS ridge group sparse* \hat{f} défini par (1.13) avec $r_v = 1$ dans (1.12). Sous les hypothèses H1, H2 et H3, il existe une constante positive C et $0 < \eta < 1$ (η tend vers 0 lorsque n augmente) tels que,

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\}, \quad (1.22)$$

avec une probabilité supérieure à $1 - \eta$.

Commentons ce résultat 1:

- R1 Soit f' la fonction dans \mathcal{F} pour laquelle l'infimum du membre de droite de l'inégalité (1.22) est réalisé. Le terme $\|m - f'\|_n^2$ est le terme de biais habituel. Il quantifie à la fois les propriétés d'approximation du RKHS \mathcal{H} et le compromis biais-variance.
- R2 Ce résultat est similaire à celui obtenu dans le cas où ε est gaussienne mais avec l'hypothèse supplémentaire (1.21). Cette hypothèse permet d'obtenir la même vitesse de convergence pour l'estimateur RKHS ridge group sparse que dans le cas où ε est gaussienne (voir [Huet and Taupin \(2017\)](#)). Cependant, elle implique certaines restrictions sur la régularité du RKHS \mathcal{H} . En effet, comme pour tout $v \in \mathcal{P}$, $\lambda_{n,v} \geq \nu_{n,v}$ (voir l'équation (1.19)), il s'ensuit que $\sum_{v \in S_f} \nu_{n,v}^2 \leq C_3 n^{2\beta-1}$, ce qui implique certaines restrictions sur la régularité du RKHS: si β est petit, ce qui sera le cas si α est grand, alors le RKHS devra être très régulier.
- R3 Par l'équation (1.19), on a aussi que pour tout $v \in \mathcal{P}$, $\lambda_{n,v} \geq \sqrt{d/n}$. Cette hypothèse permet de contrôler la probabilité de $|\mathcal{P}|$ événements (voir l'équation (3.48) du chapitre 3), où $\log(|\mathcal{P}|)$ est d'ordre d .
- R4 Le résultat 1 peut être généralisé au cas où $\sigma \neq 1$ dans l'équation (1.1), et où $r_v \neq 1$ dans l'équation (1.12), voir la remarque 3.3.5 du chapitre 3 pour une brève démonstration de ce point.

Résultat 2: la majoration du risque quadratique de l'estimateur *RKHS ridge group sparse*

Sous les mêmes hypothèses que pour le résultat 1, on a avec une grande probabilité pour une constante positive C' que,

$$\|m - \hat{f}\|_2^2 \leq C' \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \|m - f\|_2^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\}.$$

Remark 1.2.2 Le résultat 2 peut être généralisé au cas où $\sigma \neq 1$ dans l'équation (1.1), et où $r_v \neq 1$ dans l'équation (1.12) (voir la remarque 3.3.6 du chapitre 3 pour plus de détails sur ce point).

Vitesse de convergence

Sous les mêmes hypothèses que pour le résultat 1, on a:

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \left(\sum_{v \in S_f} \nu_{n,v}^2 + \frac{d|S_f|}{n} \right) \right\}.$$

Cette inégalité met en évidence que la borne supérieure est pertinente lorsque l'infimum est atteint pour les fonctions f qui ont une décomposition sparse dans \mathcal{H} , c'est-à-dire $|S_f|$ est petit, et lorsque d est petit devant n . Lorsque d est grand,

la décomposition des fonctions dans \mathcal{H} doit être limitée aux interactions d'un ordre limité, de sorte que le nombre d'éléments dans le méta-modèle estimé soit d'un ordre inférieur à d^r pour un petit r , disons $r = 2$ par exemple. Dans ce cas, le cardinal de \mathcal{P} sera donc inférieur à d^2 . Comme indiqué dans la remarque R3, l'hypothèse $\lambda_{n,v} \geq \sqrt{d/n}$ est nécessaire pour contrôler la valeur $\log(|\mathcal{P}|)$, qui sera désormais inférieur à $2 \log(d)$. Par conséquent, la valeur d dans la définition de $\lambda_{n,v}$ (voir l'équation (1.19)) ainsi que le terme $d|S_f|/n$ dans l'infimum ci-dessus seront remplacés par $2 \log(d)$ et $2 \log(d)|S_f|/n$, respectivement.

1.2.4 Travaux antérieurs

Plusieurs auteurs ont étudié les propriétés théoriques d'estimateurs similaires à l'estimateur *RKHS ridge group sparse*. Rappelons brièvement leur cadre de travail et leurs résultats.

Meier et al. (2009) ont considéré un estimateur similaire à l'estimateur *RKHS ridge group sparse*. Au lieu d'ajouter deux pénalités distinctes de sparsité et de régularité, ils combinent les deux termes en une seule pénalité de sparsité et de régularité. Ils considèrent un modèle de régression où les variables X_1, \dots, X_d sont contrôlées (non aléatoires) et où l'erreur ε est de distribution sous-gaussienne. Ils ont établi des majorations du risque empirique pour l'estimation de m sur l'ensemble des fonctions additives univariées. Par la suite, Raskutti et al. (2012) ont montré (dans la section 3.4. de leur article) que la vitesse de convergence de cet estimateur est sous-optimale.

Koltchinskii and Yuan (2010) ont considéré un estimateur de type *ridge group sparse* défini sur un ensemble de fonctions additives dont chaque terme appartient à un espace RKHS. Ils ne supposent pas que les variables d'entrées X_1, \dots, X_d sont indépendantes, ni l'orthogonalité entre les espaces RKHS. Par contre, ils introduisent des hypothèses liées au degré de dépendance des RKHS, ce qui assurant ainsi une *quasi orthogonalité* entre ces espaces. Sous l'hypothèse où $\sup_{f \in \mathcal{H}} \sup_{X \in \mathcal{X}} |f(X)|$ est borné indépendamment de la dimension d , ils ont établi des majorations de l'excès du risque en supposant que la fonction m a une représentation sparse. Leurs résultats sont valables pour une grande classe de fonctions de perte, appelée *pertes de type quadratique* qui doivent satisfaire des conditions de bornitude sur le support de la variable de sortie Y . La section 2.1. de leur article donne plusieurs exemples du cadre d'application de leurs résultats. Il convient de noter que la fonction de perte quadratique dans le cas où Y n'est pas bornée ne satisfait pas les conditions dans Koltchinskii and Yuan (2010). Les preuves de leurs résultats reposent sur les résultats établis pour la symétrisation et les inégalités de concentration pour les processus de Rademacher et sur les bornes exponentielles de type Bernstein.

Raskutti et al. (2012) ont supposé que la fonction m a une représentation additive univariée et sparse (telle que définie dans l'équation (1.9)) de sorte que chaque fonction univariée se trouve dans un RKHS. Ils ont proposé la procédure *ridge group sparse* pour calculer l'estimateur de m , et ont étudié les propriétés théoriques de leur estimateur dans le modèle de régression gaussienne. Ils ont fourni les majora-

tions du risque empirique L^2 et du risque quadratique et une minoration du risque quadratique de leur estimateur sur des espaces de modèles additifs sparse, y compris les polynômes, les splines et les classes de Sobolev.

Huet and Taupin (2017) ont étudié les propriétés théoriques de l'estimateur *RKHS ridge group sparse*, dans le modèle de régression gaussienne. Elles ont établi les majorations du risque empirique L^2 et du risque quadratique de l'estimateur *RKHS ridge group sparse*, c'est-à-dire des bornes supérieures par rapport à la norme L^2 et à la norme empirique L^2 pour la distance entre la fonction m et son estimation dans le RKHS \mathcal{H} .

Raskutti et al. (2012) et Huet and Taupin (2017) ne supposent pas que la quantité $\sup_{f \in \mathcal{H}} \sup_{X \in \mathcal{X}} |f(X)|$ est borné. Par contre, ils considèrent l'hypothèse H1 où pour tout $v \in \mathcal{P}$, $\sup_{X_v} |f_v(X_v)|$ est borné. Les preuves de leurs résultats reposent sur les méthodes probabilistes des processus empiriques gaussiens telles que les inégalités de concentration et la minoration de Sudakov (Pisier (1989), Massart (2000), van de Geer et al. (2000), Ledoux (2001)), ainsi que sur les résultats sur la complexité Rademacher des classes de noyau (Mendelson (2002), Bartlett et al. (2005)).

1.2.5 Outils techniques pour les preuves

Dans ce travail, les majorations du risque empirique L^2 et du risque quadratique de l'estimateur *RKHS ridge group sparse* sont fournies, dans le modèle de régression où l'erreur ε est non-gaussienne et non-bornée, et en considérant un critère des moindres carrés pénalisé. Dans ce cas les conditions de Koltchinskii and Yuan (2010) et les méthodes probabilistes habituelles des processus empiriques gaussiens telles que les inégalités de concentration et la minoration de Sudakov ne s'appliquent pas. Les preuves de nos résultats nécessitent des outils mathématiques différents de ceux utilisés dans les travaux précédents:

- ✓ une minoration de type Sudakov pour des variables aléatoires non-gaussiennes et non-bornées,
- ✓ une inégalité de concentration pour les queues inférieures et supérieures d'une fonction convexe des variables aléatoires non-gaussiennes et non-bornées.

A notre connaissance, dans notre contexte non-gaussien et non-borné, et avec le critère des moindres carrés, la seule minoration de type Sudakov qui permette d'obtenir la même vitesse de convergence pour l'estimateur *RKHS ridge group sparse* que dans le modèle de régression gaussienne (voir Huet and Taupin (2017)), est celle donnée par Talagrand (1994). La minoration donnée par Talagrand (1994) est spécifique aux densités π_α (voir l'équation (1.17)). C'est la raison pour laquelle la classe de densité \mathcal{D} est considérée dans ce travail. La minoration de type Sudakov adaptée à notre travail est déduite du théorème 3.1. de Talagrand (1994). Rappelons cette minoration.

Soient $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ des variables aléatoires i.i.d. distribuées avec la densité $\pi_\alpha \in \mathcal{D}$, et pour une fonction $g : \mathbb{R}^{|v|} \mapsto \mathbb{R}$, $v \in \mathcal{P}$ appartenant à une classe de

fonctions \mathcal{G} , soit $V_{n,\varepsilon}$ le processus empirique associé au vecteur aléatoire ε ,

$$V_{n,\varepsilon}(g) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i g(X_{v,i}). \quad (1.23)$$

Alors, pour tout $\delta > 0$,

$$\begin{aligned} \frac{1}{K} \log N(\delta, \mathcal{G}, \|\cdot\|) &\leq \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^2 \mathbf{1}_{[2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|, \infty)}(\delta) \\ &+ \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^\alpha \mathbf{1}_{(0, 2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|]}(\delta), \end{aligned} \quad (1.24)$$

où K est une constante qui ne dépend que de α , $N(\delta, \mathcal{G}, \|\cdot\|)$ est le nombre de recouvrements de l'espace métrique $(\mathcal{G}, \|\cdot\|)$ par des boules de rayon inférieur à δ , et $\mathbf{1}_A : \mathcal{A} \rightarrow \{0, 1\}$ est la fonction caractéristique de $A \subset \mathcal{A}$,

$$\mathbf{1}_A(a) = \begin{cases} 1 & \text{if } a \in A, \\ 0 & \text{if } a \notin A. \end{cases}$$

Concernant l'inégalité de concentration, il est démontré que les fonctions de distribution associées aux densités $\pi_\alpha \in \mathcal{D}$ appartiennent à une classe de fonctions de distribution définie par [Adamczak \(2005\)](#), pour laquelle l'inégalité de *log-Sobolev* ([Gross \(1975\)](#)) est satisfaite (voir lemme 3.4.2 du chapitre 3). [Shu and Strzelecki \(2017\)](#) ont établi des bornes pour les queues inférieures et supérieures de fonctions convexes de variables aléatoires indépendantes qui satisfont l'inégalité de *log-Sobolev*. Comme les fonctions de distribution associées aux densités $\pi_\alpha \in \mathcal{D}$ satisfont l'inégalité de *log-Sobolev*, l'inégalité de concentration dérivée par [Shu and Strzelecki \(2017\)](#) est valable pour elles. L'inégalité de concentration adaptée à notre travail est déduite du corollaire 1.7. de [Shu and Strzelecki \(2017\)](#) (voir corollaire 3.4.2 du chapitre 3).

1.3 Résumé du chapitre 4

1.3.1 Objectifs et résultats

Un package R, appelé **RKHSMetaMod**, a été développé pour mettre en œuvre la procédure *ridge group sparse* décrite dans la section 1.1.5.1. Ce package permet de:

- ✓ calculer les noyaux auto-reproduisants comme décrit dans la section 1.1.4.2, et leurs matrices de Gram associées,
- ✓ mettre en œuvre la procédure *RKHS ridge group sparse* et un cas particulier de celle-ci appelé *RKHS group lasso* (lorsque $\gamma_v = 0$, $v \in \mathcal{P}$ dans le critère (1.14)) afin d'estimer les termes f_v^* dans la décomposition de Hoeffding de f^* conduisant à une estimation de la fonction m ,
- ✓ choisir les paramètres de régularisation μ_v, γ_v , $v \in \mathcal{P}$ dans le critère (1.14) en utilisant une procédure qui permet d'obtenir le *meilleur* estimateur *RKHS ridge group sparse* en termes de qualité de prédiction,

- ✓ d'estimer les indices de Sobol de la fonction m comme décrit à la section 1.1.5.2.

Le package **RKHSMetaMod** fournit une interface entre l'environnement de calcul statistique R et les bibliothèques C++ **Eigen** et **GSL**. Afin d'optimiser le temps de calcul et la mémoire de stockage, toutes les fonctions de ce package ont été écrites en utilisant les bibliothèques **Eigen** et **GSL** de C++ à l'exception d'une fonction qui est écrite en R. Elles sont ensuite interfacées avec l'environnement R afin de proposer un package facilement exploitable aux utilisateurs de R. Le package **RKHSMetaMod** est dédié à l'estimation du méta-modèle f^* d'un modèle m sur le RKHS \mathcal{H} . Les algorithmes d'optimisation convexe utilisés dans ce package sont adaptés pour prendre en compte le problème de la grande dimensionnalité dans ce contexte. Ce package est disponible sur le Comprehensive R Archive Network (CRAN) à <https://cran.r-project.org/web/packages/RKHSMetaMod/>.

1.3.2 Présentation du modèle

Considérons un phénomène décrit par un modèle m dépendant de d variables d'entrée $X = (X_1, \dots, X_d)$. Ce modèle m de \mathbb{R}^d vers \mathbb{R} peut être un modèle connu qui est calculable en tout point X , ou un modèle de regression comme défini dans l'équation (1.1). Dans le second cas, l'erreur ε est supposée être centrée avec une variance finie. Les coordonnées de X sont indépendantes et ont la loi uniforme sur $\mathcal{X} = [0, 1]^d$. C'est-à-dire $X \sim P_X = \bigotimes_{a=1}^d P_a$, où chaque P_a , $a = 1, \dots, d$ représente la loi uniforme sur l'intervalle $[0, 1]$. Le modèle m peut être complexe, présenter de fortes non-linéarités et des effets d'interaction d'ordre élevé, et il est supposé être de carré-intégrable sur \mathcal{X} .

1.3.3 Critère à minimiser

Considérons la forme paramétrique du critère *RKHS ridge group sparse* défini dans l'équation (1.14), où γ_v et μ_v , $v \in \mathcal{P}$ sont choisis comme suit:

Pour chaque $v \in \mathcal{P}$, soient γ'_v et μ'_v les poids qui sont choisis de manière appropriée. Alors,

$$\gamma_v = \gamma \times \gamma'_v \text{ et } \mu_v = \mu \times \mu'_v, \gamma, \mu \in \mathbb{R}^+.$$

Remark 1.3.1 *Cette formulation simplifie le choix des paramètres de régularisation, car au lieu de sélectionner les paramètres γ_v et μ_v pour tous les $v \in \mathcal{P}$, seuls deux paramètres γ et μ sont sélectionnés. De plus, les poids γ'_v et μ'_v , $v \in \mathcal{P}$, peuvent être intéressants pour les applications. Par exemple, on peut prendre des poids qui augmentent avec le cardinal de v afin de favoriser les effets avec un ordre d'interaction petit entre les variables.*

Pour plus de simplicité, dans le reste de ce chapitre les valeurs de γ'_v et μ'_v pour tout $v \in \mathcal{P}$ sont fixées à 1, et le critère *RKHS ridge group sparse* est alors exprimé comme suit:

$$C(f_0, \theta) = \left\{ \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \gamma \sum_{v \in \mathcal{P}} \|K_v \theta_v\| + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\| \right\}.$$

En ne considérant que la deuxième partie de la fonction de pénalité dans le critère ci-dessus, c'est-à-dire en fixant γ à zéro, on obtient le critère du *RKHS group lasso* comme suit,

$$C_g(f_0, \theta) = \left\{ \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\| \right\}.$$

À une transformation près, $\beta_v = K_v^{1/2} \theta_v$, le critère C_g est exactement un critère de *group lasso* (Yuan and Lin (2006)).

Il convient de préciser que, dans le package **RKHSMetaMod**, les solutions de l'algorithme du *RKHS group lasso* sont utilisées afin d'initialiser les paramètres d'entrées de l'algorithme du *RKHS ridge group sparse*. En effet, la fonction de pénalité dans le critère de *RKHS group lasso* $C_g(f_0, \theta)$ assure la sparsité de la solution. Ainsi, pour une valeur donnée de μ , en implémentant l'algorithme du *RKHS group lasso*, on obtient une solution avec peu de termes dans sa décomposition additive.

Le paramètre de régularisation dans l'algorithme du *RKHS group lasso* sera noté par:

$$\mu_g = \sqrt{n} \mu. \quad (1.25)$$

1.3.3.1 Choix des paramètres de régularisation

Lorsque on est confronté à un problème d'optimisation, l'une des étapes essentielles consiste à choisir correctement les paramètres de régularisation. Pour cela,

- ✓ d'abord une grille de valeurs de μ et γ est choisie.

Soit μ_{\max} la valeur la plus petite de μ_g (voir équation (1.25)), de sorte que la solution à la minimisation du problème de *RKHS group lasso* pour tout $v \in \mathcal{P}$ est $\theta_v = 0$. On a,

$$\mu_{\max} = \max_v \left(\frac{2}{\sqrt{n}} \|K_v^{1/2} (Y - \bar{Y})\| \right).$$

Pour configurer la grille de valeurs de μ , il suffit de trouver μ_{\max} , puis une grille de valeurs de μ est définie comme suit:

$$\mu_l = \frac{\mu_{\max}}{(\sqrt{n} \times 2^l)}, \quad l \in \{1, \dots, l_{\max}\}.$$

La grille de valeurs de γ est choisie par l'utilisateur.

- ✓ ensuite, pour la grille de valeurs de μ et γ , une suite d'estimateurs est calculée. Chaque estimateur associé à la paire (μ, γ) dans la grille de valeurs de μ et γ , noté par $\hat{f}_{(\mu, \gamma)}$, est la solution du problème d'optimisation de *RKHS ridge group sparse* ou du problème d'optimisation de *RKHS group lasso* si $\gamma = 0$.
- ✓ enfin, les estimateurs $\hat{f}_{(\mu, \gamma)}$ sont évalués à l'aide d'un ensemble de données de test,

$$\{(Y_i^{\text{test}}, X_i^{\text{test}})\}_{i=1}^{n^{\text{test}}}.$$

L'erreur de prédiction associée à l'estimateur $\hat{f}_{(\mu, \gamma)}$ est calculée par,

$$\text{ErrPred}(\mu, \gamma) = \frac{1}{n^{\text{test}}} \sum_{i=1}^{n^{\text{test}}} (Y_i^{\text{test}} - \hat{f}_{(\mu, \gamma)}(X_i^{\text{test}}))^2,$$

où pour $S_{\hat{f}}$ étant le support de l'estimateur $\hat{f}_{(\mu, \gamma)}$,

$$\hat{f}_{(\mu, \gamma)}(X^{\text{test}}) = \hat{f}_0 + \sum_{v \in S_{\hat{f}}} \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v^{\text{test}}).$$

La paire $(\hat{\mu}, \hat{\gamma})$ avec la plus petite valeur de l'erreur de prédiction est choisie, et l'estimateur $\hat{f}_{(\hat{\mu}, \hat{\gamma})}$ est considéré comme le *meilleur* estimateur de la fonction m , par rapport à l'erreur de prédiction.

Dans le package **RKHSMetaMod**, les algorithmes pour calculer une suite d'estimateurs \hat{f} , la valeur de μ_{\max} et l'erreur de prédiction sont implémentés respectivement dans les fonctions `RKHSMetMod`, `mu_max` et `PredErr`.

1.3.3.2 Estimation des indices de Sobol

Les indices de Sobol de la fonction m sont estimés par les indices de Sobol empiriques de l'estimateur \hat{f} comme décrit dans la section 1.1.5.2,

$$\hat{S}_v = \begin{cases} \frac{\widehat{\text{var}}(\hat{f}_v(X_v))}{\sum_{v \in \mathcal{P}} \widehat{\text{var}}(\hat{f}_v(X_v))} & \text{pour } v \in S_{\hat{f}}, \\ 0 & \text{pour } v \notin S_{\hat{f}}. \end{cases}$$

Dans le package **RKHSMetaMod**, l'algorithme permettant de calculer des indices empiriques de Sobol \hat{S}_v , $v \in \mathcal{P}$ est implémenté dans la fonction `SI_emp`.

1.3.4 Algorithmes

Le package **RKHSMetaMod** met en œuvre deux algorithmes d'optimisation: le *RKHS ridge group sparse* et le *RKHS group lasso*. Ces algorithmes reposent sur les matrices de Gram K_v , $v \in \mathcal{P}$ qui doivent être définies positives. Ainsi, la première étape essentielle du package **RKHSMetaMod** consiste à calculer ces matrices et à s'assurer qu'elles sont définies positives.

La deuxième étape consiste à calculer l'estimateur \hat{f} . Dans le package **RKHSMetaMod**, deux objectifs différents basés sur des procédures différentes sont considérés afin de calculer cet estimateur:

- ✓ L'estimateur ayant la *meilleure* qualité de prédiction:

Dans ce cas, le *meilleur* estimateur est calculé en utilisant la procédure décrite à la section 1.3.3.1.

- ✓ L'estimateur avec pour maximum $qmax$ groupes actifs:

Le paramètre de régularisation γ est fixé à zéro. Une valeur de μ pour laquelle le nombre de groupes dans la solution du problème du *RKHS group lasso* est égal à $qmax$ est calculée. Cette valeur sera notée par μ_{qmax} . Ensuite, l'algorithme du *RKHS ridge group sparse* est implémenté pour une grille de valeurs de $\gamma \neq 0$ et la valeur de μ_{qmax} .

Cette procédure est implémentée dans le package **RKHSMetaMod** dans la fonction `RKHSMetMod_qmax`.

1.3.4.1 Calcul des matrices de Gram

Les noyaux disponibles dans le package **RKHSMetaMod** sont: noyau linéaire, noyau quadratique, noyau brownien, noyau matérn et noyau gaussien. Le choix du noyau, par l'utilisateur, détermine l'espace d'approximation fonctionnelle. Pour un noyau choisi, l'algorithme de calcul des matrices de Gram K_v , $v \in \mathcal{P}$ dans le package **RKHSMetaMod** est implémenté dans la fonction `calc_Kv`, et est basé sur trois points essentiels:

- ✓ Modifier le noyau choisi:

Afin de satisfaire les conditions de construction du RKHS \mathcal{H} décrites dans la section 1.1.4.2, ces noyaux sont modifiés selon l'équation (1.8). Ci-dessous l'exemple du noyau brownien.

Exemple 1.3.1 *La présentation habituelle du noyau brownien est la suivante:*

$$k_a(X_a, X'_a) = \min(X_a, X'_a) + 1.$$

Le RKHS associé au noyau k_a est l'ensemble,

$$\mathcal{H}_a = \left\{ f : [0, 1] \rightarrow \mathbb{R} \text{ est absolument continu et } f(0) = 0, \int_0^1 f'(X_a)^2 dX_a < \infty \right\},$$

avec le produit scalaire

$$\langle f, h \rangle_{\mathcal{H}_a} = \int_0^1 f'(X_a)h'(X_a)dX_a.$$

Le noyau k_{0a} associé au noyau brownien est calculé comme suit,

$$\begin{aligned} k_{0a} &= \min(X_a, X'_a) + 1 - \frac{(\int_0^1 (\min(X_a, U) + 1)dU)(\int_0^1 (\min(X'_a, U) + 1)dU)}{(\int_0^1 \int_0^1 (\min(U, V) + 1)dU dV)}, \\ &= \min(X_a, X'_a) + 1 - \frac{3}{4}(1 + X_a - \frac{X_a^2}{2})(1 + X'_a - \frac{X'^2_a}{2}). \end{aligned}$$

Le RKHS associé au noyau k_{0a} est l'ensemble,

$$\mathcal{H}_{0a} = \left\{ f \in \mathcal{H}_a : \int_0^1 f(X_a) dX_a = 0 \right\}.$$

Enfin, le RKHS $\mathcal{H} = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v$ est l'ensemble suivant,

$$\mathcal{H} = \left\{ f : [0, 1]^d \rightarrow \mathbb{R} : f = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \text{ avec } f_v \in \mathcal{H}_v \right\}.$$

- ✓ Calculer les matrices de Gram K_v pour tout v :

Tout d'abord, pour tout $a = 1, \dots, d$ les matrices de Gram K_a associées aux noyaux k_{0a} sont calculées en utilisant l'équation (1.8),

$$(K_a)_{i,i'} = k_{0a}(X_{ai}, X_{ai'}).$$

Ensuite, pour tout $v \in \mathcal{P}$, les matrices de Gram K_v associées au noyau $k_v = \prod_{a \in v} k_{0a}$ sont calculées comme suit:

$$K_v = \bigodot_{a \in v} K_a,$$

où \bigodot dénote le produit matriciel de Hadamard.

- ✓ Assurer que les matrices K_v , $v \in \mathcal{P}$ sont définies positives:

La sortie de la fonction `calc_Kv` est l'un des paramètres d'entrées des fonctions associées aux algorithmes *RKHS group lasso* et *RKHS ridge group sparse*. Comme ces deux algorithmes reposent sur la positivité de ces matrices, il est indispensable que les matrices K_v , $v \in \mathcal{P}$ soient définies positives. Pour cela, la fonction `calc_Kv` modifie les valeurs propres de la matrice K_v , $v \in \mathcal{P}$ si nécessaire.

Pour chaque groupe $v \in \mathcal{P}$, soit $\lambda_{v,\max}$ et $\lambda_{v,\min}$ respectivement le maximum et le minimum des valeurs propres associées à la matrice K_v , et soit "tol" un scalaire positif à fixer. Pour chaque matrice K_v ,

$$\text{"si } \lambda_{v,\min} < \lambda_{v,\max} \times \text{tol"}$$

alors, $\lambda_{v,\max} \times \text{tol}$ est ajouté à toutes les valeurs propres de K_v .

La valeur de "tol" est fixée par défaut à $1e^{-8}$, mais on peut considérer une valeur plus petite ou plus grande en fonction du noyau choisi et de n .

1.3.4.2 Algorithmes d'optimisation

RKHS group lasso Afin de résoudre le problème d'optimisation du *RKHS group lasso*, l'algorithme classique de *block coordinate descent* est utilisé (Boyd et al. (2011), Bubeck (2015)). La minimisation du critère $C_g(f_0, \theta)$ se fait à travers chaque

groupe v à chaque fois. À chaque étape de l'algorithme, le critère est minimisé en fonction des paramètres du bloc actuel, tandis que les valeurs des paramètres des autres blocs sont fixées à leurs valeurs actuelles. La procédure est répétée jusqu'à convergence.

Dans le package **RKHSMetaMod**, l'algorithme classique de *block coordinate descent* pour résoudre le problème d'optimisation du *RKHS group lasso* est implémenté dans la fonction `RKHSgrplasso`.

RKHS ridge group sparse Afin de résoudre le problème d'optimisation *RKHS ridge group sparse*, un algorithme adapté de *block coordinate descent* est proposé. Cet algorithme fournit deux étapes:

Étape 1 Initialiser les paramètres d'entrées par les solutions de l'algorithme *RKHS group lasso* pour chaque valeur du paramètre de régularisation μ et exécuter l'algorithme *RKHS ridge group sparse* via le support actif des solutions *RKHS group lasso* jusqu'à ce qu'il atteigne la convergence.

Cette étape est prévue afin de diminuer le temps de calcul.

Étape 2 Réinitialiser les paramètres d'entrées avec les solutions obtenues à l'Étape 1 et implémenter l'algorithme *RKHS ridge group sparse* à travers tous les groupes de \mathcal{P} jusqu'à ce qu'il atteigne la convergence.

Cette deuxième étape permet de vérifier qu'aucun groupe ne manque dans la sortie de l'Étape 1.

L'algorithme adapté de *block coordinate descent* pour résoudre le problème d'optimisation *RKHS ridge group sparse* est implémenté dans le package **RKHSMetaMod**, dans la fonction `pen_MetMod`.

1.4 Résumé et perspectives

Les travaux présentés dans cette thèse portent sur le problème de l'estimation d'un méta-modèle qui approche la décomposition de Hoeffding d'un modèle complexe, noté m . Le modèle m dépend de d variables d'entrée X_1, \dots, X_d qui sont indépendantes et ont une loi connue. Le méta-modèle appartient à un RKHS \mathcal{H} , qui est construit de telle manière que la décomposition additive de toute fonction f dans \mathcal{H} est la décomposition de Hoeffding de f (Durrande et al. (2013)). L'estimateur du méta-modèle, noté \hat{f} , minimise un critère des moindres carrés pénalisé par une fonction de pénalité qui est la somme de la norme de Hilbert et de la norme empirique L^2 . Cette procédure, appelée *RKHS ridge group sparse*, permet à la fois de sélectionner et d'estimer les termes importants de la décomposition de Hoeffding du méta-modèle, et donc de sélectionner les indices de Sobol non-nuls et de les estimer (Huet and Taupin (2017)).

La première partie de ce travail est dédiée à l'étude des propriétés théoriques de l'estimateur \hat{f} d'un modèle de régression où l'erreur ε est non-gaussienne et non-

bornée. Les majorations du risque empirique L^2 et du risque quadratique de cet estimateur sont fournies.

Dans la deuxième partie de ce travail, la procédure de calcul de \hat{f} est mise en œuvre dans un package R, appelé **RKHSMetaMod**. Afin d’optimiser le temps de calcul et la mémoire de stockage, toutes les fonctions de ce package ont été écrites en utilisant les bibliothèques **GSL** et **Eigen** de C++ à l’exception d’une fonction qui est écrite en R. Elles sont ensuite interfacées avec l’environnement R afin de proposer un package facilement exploitable aux utilisateurs de R. Le package **RKHSMetaMod** s’applique indifféremment dans le cas où le modèle m est calculable et le cas du modèle de régression. Une étude de simulation est fournie afin de valider la performance des fonctions du package en termes de qualité prédictive de l’estimateur et d’estimation des indices de Sobol.

Comme tous les travaux de recherche qui sont menés dans un temps limité, de nombreuses pistes n’ont pas été explorées dans ce travail et il y a plusieurs perspectives à considérer pour une étude plus approfondie. Mentionnons en quelques-unes.

1.4.1 Variables d’entrée non-indépendantes

Dans les deux parties de cette thèse, les variables d’entrée X_1, \dots, X_d du modèle m sont supposées indépendantes et leur loi est connue. Sous ces hypothèses, il est possible de construire des espaces d’approximation tels que toute fonction dans ces espaces se décompose selon sa décomposition de Hoeffding. La décomposition est unique et les termes de cette décomposition sont orthogonaux.

Si les variables X_1, \dots, X_d ne sont pas indépendantes, il n’y a plus d’orthogonalité entre les termes de la décomposition sur les espaces d’approximation et la décomposition d’une fonction sur ces espaces n’est pas nécessairement unique. Il s’ensuit que la décomposition de la variance donnée à l’équation (1.6) n’est plus valable, ni le calcul des indices de Sobol. Néanmoins, l’approximation du modèle sur un espace fonctionnel selon une décomposition additive peut s’avérer intéressante en pratique, l’estimation du méta-modèle pouvant aider à l’interprétation des effets des variables d’entrée sur la variable de sortie.

Le cas où les variables X_1, \dots, X_d ne sont pas indépendantes a été considéré par [Koltchinskii and Yuan \(2010\)](#). Leur espace d’approximation \mathcal{H} est l’espace linéaire engendré (ou *linear span* (*l.s.*)) par un *dictionnaire* d’espaces RKHS $\mathcal{H}_1, \dots, \mathcal{H}_N$,

$$\mathcal{H} = l.s. \bigcup_{j=1}^N \mathcal{H}_j.$$

L’espace \mathcal{H} est ainsi constitué de toutes les fonctions f qui ont une représentation additive de la forme,

$$f = \sum_{j=1}^N f_j(X), \quad f_j \in \mathcal{H}_j, \quad j = 1, \dots, N. \quad (1.26)$$

Sous des hypothèses sur le degré de *dépendance* des espaces RKHS \mathcal{H}_j assurant une *quasi* orthogonalité entre ces espaces, Koltchinskii and Yuan (2010) ont établi des majorations de l'excès de risque d'un estimateur de type *ridge group sparse*.

Remarquons que les espaces d'approximation considérés dans cette thèse sont un cas particulier des espaces considérés par Koltchinskii and Yuan (2010). Cependant, le contexte dans lequel leurs résultats sont établis diffère du notre en plusieurs points. D'une part, les fonctions dans l'espace \mathcal{H} sont supposées uniformément bornées au sens suivant: $\sup_{f \in \mathcal{H}} \sup_{X \in \mathcal{X}} |f(X)|$ est borné indépendamment de la dimension d . D'autre part le modèle statistique est différent et en particulier le cas du modèle de régression à erreur additive non-bornée n'est pas considéré dans leurs travaux.

Mon objectif serait donc d'établir une borne de risque pour un estimateur *ridge group sparse* sur des espaces d'approximation construits comme proposé par Durande et al. (2013), dans un cadre où les variables d'entrée X_1, \dots, X_d ne sont pas indépendantes, pour le modèle de régression à erreur additive non-bornée. L'une des étapes essentielle de la preuve repose sur la majoration de la norme L^2 dans \mathcal{H} par la norme empirique L^2 dans \mathcal{H} (Lemme 3.5.4). Le calcul de cette majoration nécessite le contrôle des moments d'ordre 4 des fonctions de \mathcal{H} . Dans le cas où les X_1, \dots, X_d sont indépendantes, la décomposition des fonctions dans \mathcal{H} est orthogonale, et ce contrôle est aisé à obtenir. Dans le cas contraire, les hypothèses sur le degré de *dépendance* des espaces \mathcal{H}_j formulées par Koltchinskii and Yuan (2010) ne permettent pas de gérer les moments d'ordre 4. Il reste donc à poursuivre ce travail pour établir une majoration du risque de l'estimateur *RKHS ridge group sparse* sous des hypothèses qui restent encore à préciser. A notre connaissance ce cas n'a pas été étudié à ce jour.

Concernant la réalisation de l'analyse de sensibilité dans ce cas, comme le calcul des indices de Sobol n'est plus possible, on peut considérer les valeurs de Shapley (Shapley (1953)), voir par exemple Owen (2014), Song et al. (2016), Owen and Prieur (2017), Benoumechiara and Elie-Dit-Cosaque (2019), Broto et al. (2019), Iooss and Prieur (2019).

1.4.2 Généralisation au modèle de régression avec erreur log-concave

Le résultat 1 montre que la majoration du risque dans le cas où les erreurs sont de densité $\pi_\alpha \in \mathcal{D}$ est la même que celle obtenue dans le cas des erreurs gaussiennes. Cependant, la classe des densités π_α est restrictive et il serait intéressant d'obtenir un résultat pour des classes de densité plus grandes, comme les densités log-concaves par exemple.

Comme expliqué à la section 1.2.5, l'une des étapes essentielles de la preuve du résultat repose sur une minoration de type Sudakov de l'espérance du processus empirique, lorsque les variables aléatoires sont supposées non-bornées et non-gaussiennes.

Dans le cas gaussien, la minoration de Sudakov s'énonce de la façon suivante (Pisier (1989)):

Soient $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ des variables aléatoires i.i.d. gaussiennes, et pour une fonction $g : \mathbb{R}^{|v|} \mapsto \mathbb{R}$, $v \in \mathcal{P}$ appartenant à une classe de fonctions \mathcal{G} , soit $V_{n,\varepsilon}$ le processus empirique associé au vecteur aléatoire ε défini dans l'équation (1.23). Alors pour tout $\delta > 0$,

$$\frac{1}{C} \log N(\delta, \mathcal{G}, \|\cdot\|) \leq \left(\frac{nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^2, \quad (1.27)$$

où C est une constante, et $N(\delta, \mathcal{G}, \|\cdot\|)$ est le nombre de recouvrements de l'espace métrique $(\mathcal{G}, \|\cdot\|)$ par des boules de rayon inférieur à δ .

Il reste ensuite à caractériser la complexité de l'espace fonctionnel \mathcal{G} pour obtenir une minoration de l'espérance du processus empirique et en déduire le résultat dans la borne de risque.

Dans le cas où les $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ sont de densité π_α , le théorème 3.1. de Talagrand (1994) établit l'inégalité donnée à l'équation (3.34) d'où on peut déduire la minoration de l'espérance du processus empirique donnée à l'équation (1.24).

Dans le cas de ε non-gaussienne et non-bornée, une minoration de type Sudakov pour les variables aléatoires i.i.d. log-concaves est donnée par Latała (2014). Une mesure sur \mathbb{R}^n est log-concave si et seulement si elle a une densité de la forme $\exp(-\phi(x))$, où $\phi : \mathbb{R}^n \rightarrow (-\infty, \infty]$ est convexe (Borell (1974)). La minoration de type Sudakov donnée par Latała (2014) est de la forme suivante:

Soient $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ des variables aléatoires i.i.d. log-concaves, alors:

$$\frac{1}{K'} \min \left(c^2 \delta, \log N(2 \times \max(c\delta^{1/2}, c^2\delta), \mathcal{G}, \|\cdot\|) \right) \leq nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|. \quad (1.28)$$

où K' est une constante universelle, et $c = 1/\max(512c', 8)$ pour c' étant une constante universelle.

On n'a pas pu déduire de l'inégalité (1.28) la minoration de type Sudakov adaptée qui conduise à la vitesse de convergence *optimal* pour l'estimateur *RKHS ridge group sparse*. Par *optimal*, on entend la même vitesse de convergence que dans le modèle de régression gaussien (voir Huet and Taupin (2017)). C'est la raison pour laquelle, dans ce travail, les densités $\pi_\alpha \in \mathcal{D}$ sont considérées. Néanmoins, un travail supplémentaire dans cette direction ainsi qu'une recherche bibliographique mérite d'être effectués.

Introduction in english

2.1 Framework

Consider a phenomenon described by a model m depending on d input variables $X = (X_1, \dots, X_d)$. This model m from \mathbb{R}^d to \mathbb{R} , may be complex including strong non-linearities and high order interaction effects. In the classical framework of sensitivity analysis, the model m can be calculated in a finite number of points.

When the components of X are independent, the model m can be decomposed as a so-called Hoeffding decomposition. If the law of the components of X is known, this decomposition allows to perform sensitivity analysis, and more precisely to calculate the Sobol indices of m (Sobol (2001), Saltelli et al. (2009)). However, the calculation of these indices may be very difficult or even impossible, especially when the number of the input variables d is large (Iooss (2011)).

A recent approach is to approximate m by an additive meta-model involving variables X_1, \dots, X_d and interactions between them, as proposed by Durrande et al. (2013). This meta-model, denoted f^* , is the orthogonal projection of m on a reproducing kernel Hilbert space (RKHS), denoted \mathcal{H} . The space \mathcal{H} is associated with a so-called ANOVA kernel which is defined in order to obtain the analytical expression of the terms of the Hoeffding decomposition of the functions of \mathcal{H} . As f^* is the orthogonal projection of m on \mathcal{H} , each term in its decomposition is an approximation of the associated term in the Hoeffding decomposition of m .

When d , the number of the input variables is large, the total number of terms in the Hoeffding decomposition of f^* becomes very high. One solution is to calculate a sparse approximation of f^* using penalized least-squares criterion as it is done in the non-parametric regression framework.

In this thesis, two frameworks are considered: the classical framework of sensitivity analysis where $m(X)$ is calculable in all points X , and the regression framework where m is unknown and so can not be calculated.

In the second case, for a given X , the value of $m(X)$ with respect to an error term ε is observable. Therefore, we have the observations Y such that,

$$Y = m(X) + \sigma\varepsilon, \sigma > 0. \quad (2.1)$$

As in the classical framework of sensitivity analysis, the idea is to approximate the Hoeffding decomposition of m by the meta-model f^* , and then calculate a sparse estimator of f^* using non-parametric estimation approaches. This estimator, denoted \hat{f} , is the solution of a least-squares minimization problem penalized by a

penalty function that imposes sparsity and smoothness. The construction of the estimator \hat{f} allows to estimate easily the Sobol indices of m .

This thesis consists of a theoretical part and a practical part:

- In the theoretical part, I established the upper bounds of the empirical L^2 risk and the L^2 risk of the estimator \hat{f} of a regression model as described in Equation (2.1) with error ε that is non-Gaussian and non-bounded. That is, the upper bounds with respect to the L^2 -norm and the empirical L^2 -norm for the distance between the true function m and its estimation \hat{f} into the RKHS \mathcal{H} . This part is presented in Chapter 3.
- In the practical part, I developed an R package, called RKHSMetaMod, for implementing the estimation methods of the meta-model f^* of a model m . This package deals both with the case where m is calculable and the case of the regression model. This part is presented in Chapter 4 and Appendix A.
 - In Chapter 4, the estimation methods and the algorithms used in the package are described. The performances of the package functions in terms of the predictive quality of the estimator and the estimation of the Sobol indices, are validated by a simulation study.
 - In Appendix A, the complete documentation of the package, including detailed explanations of the package functions and the examples of usage of each function of the package, is provided.

The summaries of Chapters 3 and 4 are presented in Sections 2.2 and 2.3, respectively. Before that, several tools that are common to these two Chapters are briefly described. More precisely:

- introduction to the sensitivity analysis (see Section 2.1.1),
- focus on the variance based methods of global sensitivity analysis (see Section 2.1.2),
- introduction to the meta-modelling (see Section 2.1.3),
- construction of a meta-model by projection on the reproducing kernel Hilbert spaces (RKHS) (see Section 2.1.4),
- the estimation method (see Section 2.1.5).

2.1.1 Introduction to the sensitivity analysis

The sensitivity analysis methods allow to study the relationships between the output and input variables of the model, and measure the effect of each variable or groups of variables on the model output. The underlying goals for sensitivity analysis are model calibration, model validation and assisting with the decision making process. Most of the classical methods and objectives of the sensitivity analysis can be found

in [Cacuci \(2003\)](#), [Fang et al. \(2005\)](#), [Dean and Lewis \(2006\)](#), [de Rocquigny et al. \(2008\)](#), [Saltelli \(2008\)](#), [Helton \(2008\)](#), [Saltelli et al. \(2009\)](#), [Faivre et al. \(2013\)](#), [Borgonovo and Plischke \(2016\)](#).

The sensitivity analysis procedure implies the computation and analysis of some measures that evaluate the effect of the input variables on the model output. For example, the effect of an input variable on the model output can be evaluated by the amount of variance in the model output caused by that input variable. The sensitivity analysis methods can be classified in two main groups:

Local sensitivity analysis that studies the local impact of the input variables on the output variable. It consists in calculating the gradient of the output variable with respect to the input variables around a chosen value (the mean value of the input variables for example). Numerous methods have been developed to compute the gradient efficiently, including Adjoint modelling ([Cacuci \(2003\)](#), [Cacuci and Navon \(2005\)](#)) and Automated Differentiation ([Griewank and Walther \(2008\)](#)). Local methods do not fully explore the space of input variables, since they study the impact of small perturbations of input variables (generally one variable at a time) on the output variable.

Global sensitivity analysis calculates the uncertainty of the output variable due to the variations in the input variables or groups of input variables. In contrast to the local sensitivity analysis, this class of methods considers the whole variation range of the input variables. The global sensitivity analysis methods are numerous, see for example [Saltelli et al. \(2009\)](#) for a good state-of-the-art and [Iooss and Lemaître \(2015\)](#) for a review of these methods. Generally, the methods of the global sensitivity analysis which allow to calculate the most used quantitative sensitivity measures can be gathered into two groups:

- ✓ Regression-based methods are suitable when the model is linear, i.e. if the coefficient of determination R^2 is close to one. The commonly used sensitivity measures in this case are: the standardized regression coefficients, the Pearson correlation coefficients, and partial correlation coefficients. For a non-linear model that is monotonic, these coefficients could be still used to represent the output sensitivities by applying a rank transformation ([Saltelli et al. \(2009\)](#)). When the model is non-linear and non-monotonic these methods fail to produce satisfactory sensitivity measures ([Saltelli and Sobol \(1995\)](#)).
- ✓ Variance-based methods can be applied to non-linear and non-monotonic models. They consist of decomposing the variance of the model output into parts attributable to each input variable and groups of them (interactions). The sensitivity measures in this case are expressed as the ratio of the variance of each input variable or groups of them over the variance of the model output. The decomposition of variance is meaningful if the input variables are independent from one another ([Saltelli and Tarantola \(2002\)](#)). These methods

are widely used as they allow to explore whole variation range of the input variables, accounting for interactions, and non-linear non-monotonic models.

2.1.2 Global sensitivity analysis: variance-based methods

Let us consider a model m depending on d input variables $X = (X_1, \dots, X_d)$ that are independent and distributed with a known law $P_X = P_1 \otimes \dots \otimes P_d$ on $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ a subset of \mathbb{R}^d . The model m from \mathbb{R}^d to \mathbb{R} is square-integrable, i.e. $m \in L^2(\mathcal{X}, P_X)$.

In the classical framework of sensitivity analysis, where for each value of X the value of $m(X)$ can be calculated, one may use the method of Sobol (1993) to perform sensitivity analysis on m . Let us briefly recall this method.

The independency between the components of X allows to write the model m according to its Hoeffding decomposition (Sobol (1993), van der Vaart (1998)):

$$m(X) = m_0 + \sum_{a=1}^d m_a(X_a) + \sum_{a < a'} m_{a,a'}(X_a, X_{a'}) + \dots + m_{1,\dots,d}(X). \quad (2.2)$$

The terms in this decomposition are defined in terms of the conditional expected values:

$$\begin{aligned} m_0 &= E_X(m(X)), \\ m_a(X_a) &= E_X(m(X)|X_a) - m_0, \\ m_{a,a'}(X_a, X_{a'}) &= E_X(m(X)|X_a, X_{a'}) - m_a(X_a) - m_{a'}(X_{a'}) - m_0, \end{aligned}$$

and so on for interactions of order higher than two.

These terms are known as constant term, main effects, interactions of order two and so on.

Let \mathcal{P} be the set of all subsets of $\{1, \dots, d\}$ with dimension 1 to d . For all $v \in \mathcal{P}$ and $X \in \mathcal{X}$, let X_v be the vector with components X_a , $a \in v$ and m_v be the function associated with X_v in Equation (2.2). Then Equation (2.2) can be expressed as follows:

$$m(X) = m_0 + \sum_{v \in \mathcal{P}} m_v(X_v). \quad (2.3)$$

This decomposition is unique, all the terms m_v , $v \in \mathcal{P}$ are centered, and they are orthogonal with respect to $L^2(\mathcal{X}, P_X)$, i.e.

$$\forall v \in \mathcal{P}, E_X(m_v(X_v)) = 0,$$

and

$$\forall v, v' \in \mathcal{P}, v \neq v', E_X(m_v(X_v)m_{v'}(X_{v'})) = 0.$$

The function m as well as all the functions m_v in Equation (2.3) are square-integrable. As any two terms of decomposition (2.3) are orthogonal, by squaring

(2.3) and integrating it with respect to the distribution of X , a decomposition of the variance of $m(X)$ is obtained as follows:

$$\text{var}(m(X)) = \sum_{v \in \mathcal{P}} \text{var}(m_v(X_v)). \quad (2.4)$$

For any group of variables X_v , $v \in \mathcal{P}$, the Sobol indices are defined by:

$$S_v = \frac{\text{var}(m_v(X_v))}{\text{var}(m(X))}.$$

For each v , S_v expresses the fraction of variance of $m(X)$ explained by X_v .

For all $v \in \mathcal{P}$, when $|v| = 1$, the S_v are referred to as the first order indices or main effects. When $|v| = 2$, i.e. $v = \{a, a'\}$ and $a \neq a'$, they are referred to as the second order indices or the interaction indices of order two (between X_a and $X_{a'}$). And the same holds for $|v| > 2$.

The total number of the Sobol indices to be calculated is equal to $|\mathcal{P}| = 2^d - 1$, which raises exponentially with the number of the input variables d . When d is large, the evaluation of all the indices can be too computationally demanding and even not reachable. For this reason, only the indices of order not higher than two are calculated in practice. However, only first and second order indices may not provide a good information on the model sensitivities. In order to provide a better information on the model sensitivities, Homma and Saltelli (1996) proposed to calculate the first order and the total indices defined as follows:

Let $\mathcal{P}_a \subset \mathcal{P}$ be the set of all the subsets of $\{1, \dots, d\}$ including a , then

$$S_{T_a} = \sum_{v \in \mathcal{P}_a} S_v.$$

For all $a \in \{1, \dots, d\}$, S_{T_a} denotes the total effect of X_a . It expresses the fraction of variance of $m(X)$ explained by X_a alone and all the interactions of it with the other variables.

The total indices allow to rank the input variables with respect to the amount of their effect on the output variable. However, they do not provide complete information on the model sensitivities as do all the Sobol indices.

The classical computation of the Sobol indices is based on the Monte Carlo methods (see for example: Sobol (1993) for the main effect and interaction indices, and Saltelli (2002) for the main effect and total indices). These methods are very costly, since they require many thousands of model runs to get precise estimates of the Sobol indices. Thus in the case where d is large, m is complex and the calculation of the variances is numerically complicated or not possible as in the case where the model m is unknown, the methods described above are not applicable.

Another method is to approximate m by a simplified model, called a meta-model, which is much faster to evaluate and to perform sensitivity analysis on it. A meta-model provides additional information than just scalar indices. It provides the approximations of the Sobol indices of m at a lower computational cost, and also a deeper view of the input variables effects on the model output.

2.1.3 Meta-modelling

Meta-modelling consists in building a function which is computationally tractable, easy to interpret and has good prediction qualities. Let $\{m(X_i)\}_{i=1}^n$ be the outputs of n evaluations of the model m based on an experimental design $\{X_i\}_{i=1}^n$. In this context, a meta-model is an approximation of the model m which is constructed based on the experimental design $\{X_i\}_{i=1}^n$ and the outputs $\{m(X_i)\}_{i=1}^n$. There exists different approaches of meta-modelling, see [Sacks et al. \(1989\)](#), [Friedman \(1991\)](#), [Breiman \(2001\)](#), [Friedman \(2001\)](#), [Kennedy and O'Hagan \(2001\)](#), [Oakley and O'Hagan \(2004\)](#), [Storlie and Helton \(2008\)](#), [Storlie et al. \(2009\)](#), [Storlie et al. \(2011\)](#) for different examples, and [Touzani \(2011\)](#) for an overview.

In the framework of the global sensitivity analysis, one may consider a meta-model that has the additive decomposition, and that is candidate to approximate the Hoeffding decomposition of m . This meta-model allows to perform global sensitivity analysis and calculate the Sobol indices of m , even of high order. By a function that has the additive decomposition, we mean a function f from $\mathcal{X} \subset \mathbb{R}^d$ to \mathbb{R} that is defined as follows:

$$f = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \quad E_X(f_v(X_v)) = 0, \quad E_X(f_v(X_v)f_{v'}(X_{v'})) = 0, \quad \forall v, v' \in \mathcal{P}, \quad v \neq v',$$

where f_0 is a constant, and the functions f_v are supposed to belong to some functional spaces.

Among the meta-modelling approaches proposed in the literature, the decomposition based on polynomial Chaos ([Wiener \(1938\)](#), [Schoutens \(2000\)](#)) can be used to approximate the Hoeffding decomposition of m ([Sudret \(2008\)](#)).

The principle of the polynomial Chaos is to project m onto a basis of orthonormal polynomials. The Chaos representation of m is written as ([Soize and Ghanem \(2004\)](#)):

$$m(X) = \sum_{j=0}^{\infty} h_j \phi_j(X), \quad (2.5)$$

where $\{h_j\}_{j=0}^{\infty}$ are the coefficients, and $\{\phi_j\}_{j=0}^{\infty}$ are multivariate orthonormal polynomials associated with X that are determined according to the distribution of the components of X . In practice, expansion (2.5) shall be truncated for computational purposes, and the model m may be approximated by:

$$m(X) \approx \sum_{j=0}^{v_{max}} h_j \phi_j(X),$$

where v_{max} is determined using a *truncation scheme*. In this approach, the Sobol indices are obtained by summing up the squares of the suitable coefficients.

[Blatman and Sudret \(2011\)](#) proposed a method for truncating the polynomial Chaos expansion and an algorithm based on least-angle regression for selecting the terms in the expansion.

In this method, according to the distribution of the components of X a unique family of orthonormal polynomials $\{\phi_j\}_{j=0}^{\infty}$ is determined. However, this family may not be necessarily the best functional basis to approximate m well.

Another approach to construct meta-models is given by Gaussian Process (GP) modelling (Welch et al. (1992), Oakley and O’Hagan (2004), Kleijnen (2007, 2009), Marrel et al. (2009), Durrande et al. (2012), Le Gratiet et al. (2014)). The principle is to consider that the prior knowledge about the function $m(X)$, can be modelled by a GP $\mathcal{Z}(X)$ with a mean $m_{\mathcal{Z}}(X)$ and a covariance kernel $k_{\mathcal{Z}}(X, X')$. To perform sensitivity analysis from a GP model one may replace the true model $m(X)$ with the mean of the conditional GP, and deduce the Sobol indices from it. Most of the time, with GP, the Sobol indices are estimated using Monte Carlo methods.

A review on the meta-modelling based on polynomial Chaos and GP is presented in Le Gratiet et al. (2017).

Durrande et al. (2013) considered a class of functional approximation methods similar to the GP regression and obtained a meta-model that satisfies the properties of the Hoeffding decomposition. They proposed to approximate m by functions belonging to a RKHS \mathcal{H} which is constructed as a direct sum of Hilbert spaces, such that the projection of m onto \mathcal{H} is an approximation of the Hoeffding decomposition of m .

In the regression framework, when the values of $\{m(X_i)\}_{i=1}^n$ can not be calculated, one may use the projection methods on a functional basis to estimate a meta-model of m . This meta-model is estimated based on the observations $\{(X_i, Y_i)\}_{i=1}^n$ by using non-parametric estimation approaches. The estimator obtained can be used then to estimate the Sobol indices of m . Huet and Taupin (2017) considered the same approximation functional spaces as Durrande et al. (2013), and proposed an estimator of a meta-model that approximates the Hoeffding decomposition of m . They deduced from this estimated meta-model, estimators for the Sobol indices of m . This approach is presented in more details in the following.

2.1.4 Meta-models based on the reproducing kernel Hilbert spaces (RKHS)

Let us begin this Section with a brief introduction to the RKHS. The method of Durrande et al. (2013) to construct the RKHS, and the definition of the meta-model f^* that approximates the Hoeffding decomposition of m are presented in Sections 2.1.4.2 and 2.1.4.3, respectively.

2.1.4.1 Introduction to the RKHS

Let \mathcal{H} be a Hilbert space of real valued functions on a set \mathcal{X} . The space \mathcal{H} is a RKHS if for all $X \in \mathcal{X}$ the evaluation functionals

$$\begin{aligned} L_X : \mathcal{H} &\rightarrow \mathbb{R} \\ f &\rightarrow f(X), \end{aligned}$$

are continuous.

The Riesz representation Theorem ensures the existence of a unique element $k_X(\cdot)$ in \mathcal{H} verifying the following property:

$$\forall X \in \mathcal{X}, \forall f \in \mathcal{H}, f(X) = L_X(f) = \langle f, k_X \rangle_{\mathcal{H}},$$

where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ denotes the inner product in \mathcal{H} .

It follows that for all X, X' in \mathcal{X} , and $k_X(\cdot), k_{X'}(\cdot)$ in \mathcal{H} , we have,

$$k_X(X') = L_{X'}(k_X) = \langle k_X, k_{X'} \rangle_{\mathcal{H}}. \quad (2.6)$$

This allows to define the reproducing kernel of \mathcal{H} as follows:

$$\begin{aligned} k : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ (X, X') &\mapsto k_X(X'). \end{aligned}$$

The reproducing kernel $k(X, X')$ satisfies the following properties:

- It is symmetric. Indeed, by definition of $k(\cdot, \cdot)$ and thanks to the property (2.6), we have:

$$k(X, X') = k_X(X') = \langle k_X, k_{X'} \rangle_{\mathcal{H}} = k_{X'}(X) = k(X', X).$$

- For any $n \in \mathbb{N}$, $\{X_i\}_{i=1}^n \in \mathcal{X}$ and $\{c_i\}_{i=1}^n \in \mathbb{R}$, we have:

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(X_i, X_j) &= \sum_{i=1}^n \sum_{j=1}^n \langle c_i k(X_i, \cdot), c_j k(X_j, \cdot) \rangle_{\mathcal{H}}, \\ &= \left\| \sum_{i=1}^n c_i k(X_i, \cdot) \right\|_{\mathcal{H}}^2 \geq 0. \end{aligned}$$

Thus, $k(X, X')$ is positive definite.

For more background on RKHS, we refer to various standard references such as [Aronszajn \(1950\)](#), [Saitoh \(1988\)](#), and [Berlinet and Thomas-Agnan \(2003\)](#).

2.1.4.2 RKHS construction and Hoeffding decomposition

The idea is to construct an RKHS including the functions that have the additive decomposition and that are candidate to approximate the Hoeffding decomposition of m . To do so, we use the method of [Durrande et al. \(2013\)](#) that we recall briefly in the following.

Let $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ be a subset of \mathbb{R}^d . For each $a \in \{1, \dots, d\}$, we choose a RKHS \mathcal{H}_a and its associated kernel k_a defined on the set $\mathcal{X}_a \subset \mathbb{R}$ such that the two following properties are satisfied:

- (i) $k_a : \mathcal{X}_a \times \mathcal{X}_a \rightarrow \mathbb{R}$ is $P_a \otimes P_a$ measurable,
- (ii) $E_{X_a} \sqrt{k_a(X_a, X_a)} < \infty$.

The property (ii) depends on the kernel k_a , $a = 1, \dots, d$ and the distribution of X_a , $a = 1, \dots, d$. It is not very restrictive since it is satisfied, for example, for any bounded kernel.

The RKHS \mathcal{H}_a can be decomposed as a sum of two orthogonal sub-RKHS,

$$\mathcal{H}_a = \mathcal{H}_{0a} \overset{\perp}{\oplus} \mathcal{H}_{1a},$$

where \mathcal{H}_{0a} is the RKHS of zero mean functions,

$$\mathcal{H}_{0a} = \left\{ f_a \in \mathcal{H}_a : E_{X_a}(f_a(X_a)) = 0 \right\},$$

and \mathcal{H}_{1a} is the RKHS of constant functions,

$$\mathcal{H}_{1a} = \left\{ f_a \in \mathcal{H}_a : f_a(X_a) = C \right\}.$$

The kernel k_{0a} associated with the RKHS \mathcal{H}_{0a} is defined by:

$$k_{0a}(X_a, X'_a) = k_a(X_a, X'_a) - \frac{E_{U \sim P_a}(k_a(X_a, U))E_{U \sim P_a}(k_a(X'_a, U))}{E_{(U, V) \sim P_a \otimes P_a} k_a(U, V)}. \quad (2.7)$$

Let $k_v(X_v, X'_v) = \prod_{a \in v} k_{0a}(X_a, X'_a)$, then the ANOVA kernel $k(\cdot, \cdot)$ is defined as follows:

$$k(X, X') = \prod_{a=1}^d \left(1 + k_{0a}(X_a, X'_a) \right) = 1 + \sum_{v \in \mathcal{P}} k_v(X_v, X'_v).$$

For \mathcal{H}_v being the RKHS associated with the kernel k_v , the RKHS associated with the ANOVA kernel is then defined by,

$$\mathcal{H} = \prod_{a=1}^d \left(\mathbb{1} \overset{\perp}{\oplus} \mathcal{H}_{0a} \right) = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v,$$

where \perp denotes the L^2 inner product.

According to this construction, any function $f \in \mathcal{H}$ satisfies the following decomposition:

$$f(X) = \langle f, k(X, \cdot) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v),$$

which is the Hoeffding decomposition of f .

The regularity properties of the RKHS \mathcal{H} constructed as described above, depend on the set of the kernels (k_a , $a = 1, \dots, d$). This method allows to choose different approximation spaces independently of the distribution of the input variables X_1, \dots, X_d , by choosing different sets of kernels. While as mentioned earlier, in the meta-modelling approach based on polynomial Chaos expansion, according to the distribution of the input variables X_1, \dots, X_d a unique family of orthonormal polynomials $\{\phi_j\}_{j=0}^{\infty}$ is determined. Here, the distribution of the components of X occurs only for the orthogonalization of the spaces \mathcal{H}_v , $v \in \mathcal{P}$, and not in the choice of the RKHS, under the condition that properties (i) and (ii) are satisfied. This is one of the main advantages of this method compared to the method based on the truncated polynomial Chaos expansion where the smoothness of the approximation is handled only by the choice of the truncation (Blatman and Sudret (2011)).

2.1.4.3 Approximating the Hoeffding decomposition of m

Let $f^* \in \mathcal{H}$ be the orthogonal projection of m on \mathcal{H} defined by:

$$f^* = \arg \min_{f \in \mathcal{H}} \|m - f\|_2^2 = \arg \min_{f \in \mathcal{H}} E_X (m(X) - f(X))^2.$$

The function $f^* = f_0^* + \sum_{v \in \mathcal{P}} f_v^*$ is the approximation of m on the RKHS \mathcal{H} , and its Hoeffding decomposition is an approximation of the Hoeffding decomposition of m . Therefore, for each $v \in \mathcal{P}$ the function f_v^* approximates the function m_v in Equation (2.3).

The number of functions f_v^* is related to the cardinal of \mathcal{P} , equal to $2^d - 1$, that may be huge. So, the idea is to calculate a sparse estimator of f^* as an estimator of m using non-parametric approaches.

2.1.5 Estimation method

Let us consider the regression model defined in Equation (2.1),

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0.$$

The unknown function m is approximated by the meta-model f^* which is then estimated by a sparse estimator \hat{f} . This estimator \hat{f} , based on n observations $\{(Y_i, X_i)\}_{i=1}^n$, minimizes a penalized criterion. The penalty function deals both with the non-parametric nature of the problem, and the possibly large number of functions that have to be estimated.

Before describing the method to calculate \hat{f} , let us recall some methods related to the estimation in a non-parametric additive regression model.

Some authors approximate m by a function that has a sparse univariate additive decomposition of the form,

$$f(X) = f_0 + \sum_{a \in S} f_a(X_a) \quad \text{with } |S| < d, \quad (2.8)$$

where f_0 is a constant, and for all $a \in S$ the f_a are unknown smooth functions fitted from the data.

Ravikumar et al. (2009) considered a Hilbert space \mathcal{H} of functions that have univariate additive form. Their functional approximation space \mathcal{H} is constructed as a direct sum of Hilbert spaces, i.e.

$$\mathcal{H} = \bigoplus_{a=1}^d \mathcal{H}_a,$$

where for all $a \in \{1, \dots, d\}$, \mathcal{H}_a is the Hilbert subspace of $L^2(\mathcal{X}_a, P_a)$ of P_a measurable univariate functions f_a with zero mean. In order to control smoothness and to enforce sparsity in the univariate additive decomposition, they proposed the Sparse

Additive Models (SpAM) method. Their method is based on the minimization of the least-squares criterion penalized with a penalty function defined as follows,

$$\lambda \sum_{a=1}^d \sqrt{\int (f_a(X_a))^2 dX_a}, \lambda \in \mathbb{R}^+.$$

Meier et al. (2009) proposed an estimator which lies in the space of natural cubic splines. Their method is based on the minimization of the least-squares criterion penalized with a penalty function of the form,

$$\sum_{a=1}^d \sqrt{\lambda_1 \|f_a\|_n^2 + \lambda_2 \int (f_a''(X_a))^2 dX_a}, \lambda_1, \lambda_2 \in \mathbb{R}^+,$$

where $\|f_a\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_a^2(X_{ai})$.

Their penalty function consists of two parts: the first part controls the sparsity and the second part controls the smoothness.

Raskutti et al. (2012) considered several functional approximation spaces including polynomials, splines and Sobolev. Their method is based on the minimization of the least-squares criterion penalized with a penalty function of the form,

$$\gamma \|f\|_{n,1} + \mu \|f\|_{\mathcal{H},1}, \gamma, \mu \in \mathbb{R}^+, \quad (2.9)$$

where $\|f\|_{n,1} = \sum_{a=1}^d \|f_a\|_n$, and $\|f\|_{\mathcal{H},1} = \sum_{a=1}^d \|f_a\|_{\mathcal{H}_a}$.

In their penalty function, the first part controls the sparsity and the second part controls the smoothness.

Performing global sensitivity analysis on an univariate additive model leads to obtain only the first order Sobol indices, which may not provide a good information on the model sensitivities. The interactions between variables that may affect the relationship between Y and X are completely ignored in this setting.

In order to include the interaction effects, one may approximate m by a function that has a sparse multivariate additive decomposition of the form,

$$f(X) = f_0 + \sum_{v \in S} f_v(X_v) \text{ with } |S| < |\mathcal{P}|,$$

which is a generalization of the sparse univariate additive decomposition defined in Equation (2.8).

In the framework of smoothing spline ANOVA (Wahba (1990), Friedman (1991), Wahba et al. (1995)), Lin and Zhang (2006) proposed the Component Selection and Smoothing Operator (COSSO) method. Their method is based on the minimization of the least-squares criterion penalized with a penalty function that is the combination of the l_1 -norm with the Hilbert norm. The implementation of COSSO is carried out over the second-order Sobolev spaces.

Kandasamy and Yu (2016) proposed the Shrunk Additive Least Squares Approximation (SLASA) method which is based on the minimization of the least-squares criterion penalized with the sum of squared RKHS norms. Their estimator is a

v_{max} -th order multivariate additive function containing $\binom{d}{v_{max}}$ terms in its expansion. The value of v_{max} is determined using a cross validation procedure.

Huet and Taupin (2017) considered an estimator of a meta-model that approximates the Hoeffding decomposition of m defined in Equation (2.3). Their estimator is the solution of least-squares minimization penalized by the penalty function defined in Equation (2.9) adapted to the multivariate setting,

$$\gamma \|f\|_n + \mu \|f\|_{\mathcal{H}} \text{ with } \|f\|_{\mathcal{H}} = \sum_{v \in \mathcal{P}} \|f_v\|_{\mathcal{H}_v}, \text{ and } \|f\|_n = \sum_{v \in \mathcal{P}} \|f_v\|_n.$$

Their method, called ridge group sparse, estimates the groups v that are suitable for predicting the meta-model f^* , and the relationship between f_v^* and X_v for each group $v \in \mathcal{P}$. The obtained estimator, called RKHS ridge group sparse estimator, is used then to estimate the Sobol indices of m . This method makes it possible to estimate the Sobol indices for all groups in the support of the RKHS ridge group sparse estimator, including the interactions of possibly high order, a point known to be difficult in practice.

In order to obtain a sparse estimator of the meta-model f^* in this work, we use the method proposed by Huet and Taupin (2017), which allows also to obtain the estimators of the Sobol indices of m . We recall this method and the method to estimate the Sobol indices of m in the next Section and Section 2.1.5.2, respectively.

2.1.5.1 Ridge group sparse procedure and associated estimator

For all $v \in \mathcal{P}$, let X_v be the matrix of variables corresponding to the v -th group,

$$X_v = (X_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}.$$

For any $f \in \mathcal{H}$ such that $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, and for some tuning parameters $\gamma_v, \mu_v, v \in \mathcal{P}$, the ridge group sparse criterion is defined as follows:

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \left(Y_i - f_0 - \sum_{v \in \mathcal{P}} f_v(X_{vi}) \right)^2 + \sum_{v \in \mathcal{P}} \gamma_v \|f_v\|_n + \sum_{v \in \mathcal{P}} \mu_v \|f_v\|_{\mathcal{H}_v}, \quad (2.10)$$

where $\|f_v\|_n$ is the empirical L^2 -norm of f_v defined by the sample $\{X_{vi}\}_{i=1}^n$ as follows:

$$\|f_v\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_v^2(X_{vi}).$$

The penalty function in the criterion $\mathcal{L}(f)$ is the sum of the Hilbert norm and the empirical norm, which allows to select few terms in the additive decomposition of f over sets $v \in \mathcal{P}$. Moreover, the Hilbert norm favours the smoothness of the estimated $f_v, v \in \mathcal{P}$.

Let us define the set of functions,

$$\mathcal{F} = \left\{ f : f = f_0 + \sum_{v \in \mathcal{P}} f_v, \text{ with } f_v \in \mathcal{H}_v, \text{ and } \|f_v\|_{\mathcal{H}_v} \leq r_v, r_v > 0 \right\}. \quad (2.11)$$

Then, the RKHS ridge group sparse estimator of m is defined by,

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f). \quad (2.12)$$

According to the Representer Theorem (Kimeldorf and Wahba (1970)), the non-parametric functional minimization problem (2.12) is equivalent to a parametric minimization problem. Indeed, the solution of the minimization problem (2.12) belonging to the RKHS \mathcal{H} is written as $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, where for some matrix $\theta = (\theta_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}$ we have for all $v \in \mathcal{P}$,

$$f_v(\cdot) = \sum_{i=1}^n \theta_{vi} k_v(X_{vi}, \cdot).$$

Let $\|\cdot\|$ be the Euclidean norm in \mathbb{R}^n , and for each $v \in \mathcal{P}$, let K_v be the $n \times n$ Gram matrix associated with the kernel $k_v(\cdot, \cdot)$, i.e.

$$(K_v)_{i,i'} = k_v(X_{vi}, X_{vi'}).$$

Let also $K_v^{1/2}$ be the matrix that satisfies $t(K_v^{1/2})K_v^{1/2} = K_v$, and let \hat{f}_0 and $\hat{\theta}$ be the minimizers of the following penalized least-squares criterion:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \sum_{v \in \mathcal{P}} \gamma_v \|K_v \theta_v\| + n \sum_{v \in \mathcal{P}} \mu_v \|K_v^{1/2} \theta_v\|. \quad (2.13)$$

Then, the estimator \hat{f} defined in Equation (2.12) satisfies,

$$\hat{f}(X) = \hat{f}_0 + \sum_{v \in \mathcal{P}} \hat{f}_v(X_v) \text{ with } \hat{f}_v(X_v) = \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v). \quad (2.14)$$

As criterion $C(f_0, \theta)$ is convex and separable, one may calculate $\hat{\theta}$ using a block coordinate descent algorithm (Boyd et al. (2011), Bubeck (2015)).

Remark 2.1.1 *The constraint $\|f_v\|_{\mathcal{H}_v} \leq r_v$ is not taken into account in the parametric minimization problem. This constraint is crucial for theoretical properties but the value of r_v is unknown and has no practical usefulness.*

2.1.5.2 Estimation of the Sobol indices of m

The variance of the function m is estimated by the variance of the estimator \hat{f} . As the estimator \hat{f} belongs to the RKHS \mathcal{H} , it admits the Hoeffding decomposition and,

$$\text{var}(\hat{f}(X)) = \sum_{v \in \mathcal{P}} \text{var}(\hat{f}_v(X_v)),$$

where for all $v \in \mathcal{P}$,

$$\text{var}(\widehat{f}_v(X_v)) = E_X(\widehat{f}_v^2(X_v)) = \|\widehat{f}_v\|_2^2.$$

In order to reduce the computational cost in practice, one may estimate the variances of $\widehat{f}_v(X_v)$, $v \in \mathcal{P}$ by their empirical variances.

Let \widehat{f}_v be the empirical mean of $\{\widehat{f}_v(X_{vi})\}_{i=1}^n$,

$$\widehat{f}_v = \frac{1}{n} \sum_{i=1}^n \widehat{f}_v(X_{vi}),$$

then

$$\widehat{\text{var}}(\widehat{f}_v(X_v)) = \frac{1}{n-1} \sum_{i=1}^n (\widehat{f}_v(X_{vi}) - \widehat{f}_v)^2.$$

For the groups v that belong to the support of \widehat{f} , the estimators of the Sobol indices of m are defined by,

$$\widehat{S}_v = \frac{\widehat{\text{var}}(\widehat{f}_v(X_v))}{\sum_{v \in \mathcal{P}} \widehat{\text{var}}(\widehat{f}_v(X_v))},$$

and for the groups v that do not belong to the support of \widehat{f} , we have $\widehat{S}_v = 0$.

2.2 Summary of Chapter 3

2.2.1 Objectives and results

For an estimator \widehat{f} of a model m let $R(m, \widehat{f})$ be its risk. The risk $R(m, \widehat{f})$ is a measure that characterizes the precision of the estimator \widehat{f} and that can be expressed as a function of the bias and the variance of \widehat{f} :

$$R(m, \widehat{f}) = (\text{bias}(\widehat{f}))^2 + \text{var}(\widehat{f}).$$

Thus, the quality of the estimator \widehat{f} can be measured by its risk. We consider the empirical L^2 risk of the estimator \widehat{f} , i.e. when $R(m, \widehat{f}) = \|m - \widehat{f}\|_n^2$, and the L^2 risk of the estimator \widehat{f} , i.e. when $R(m, \widehat{f}) = \|m - \widehat{f}\|_2^2$.

We are interested in non-asymptotic properties of the estimator \widehat{f} , in the sense that the number of observations n is not assumed to tend to infinity. So, our results are valid for all n with a high probability. In particular, we establish the upper bounds of the risk $R(m, \widehat{f})$ of the form,

$$R(m, \widehat{f}) \leq C \inf_{f \in \mathcal{F}} \{R(m, f) + r_n(f)\}, \quad (2.15)$$

where C is a constant, and \mathcal{F} is the approximation space.

Let f' be the function in \mathcal{F} such that the infimum of the right hand side of the inequality (2.15) is realized. The term $R(m, f')$ is the bias term which depends on

the choice of the approximation space. The term $r_n(f)$ is the variance term which has to decrease with n . It gives the rate of convergence of the estimator \hat{f} , i.e. the speed at which the estimator \hat{f} approaches the true function m . The variance term depends on the regularity of the kernels k_v , $v \in \mathcal{P}$, the number of terms involved in the decomposition of the function f on the approximation space, the number of input variables d , and the number of observations n .

In the Gaussian regression framework, i.e. when ε in Equation (2.1) is distributed as a centered Gaussian random variable, [Huet and Taupin \(2017\)](#) established the upper bounds of the empirical L^2 risk and the L^2 risk of the RKHS ridge group sparse estimator.

In this Chapter, we consider the regression framework with error ε that is non-Gaussian and non-bounded. In this context, the objective is to establish the risk upper bounds of the RKHS ridge group sparse estimator of the form (2.15) with the same rate of convergence as in the Gaussian regression framework. The upper bounds of the empirical L^2 risk and the L^2 risk of the estimator \hat{f} are presented in [Result 1](#) and [Result 2](#), respectively.

2.2.2 Presentation of the model

Consider the regression model defined in Equation (2.1),

$$Y = m(X) + \sigma\varepsilon, \sigma > 0.$$

The input variables $X = (X_1, \dots, X_d)$ are independent and have a known law $P_X = \bigotimes_{a=1}^d P_a$ on $\mathcal{X} = \prod_{a=1}^d \mathcal{X}_a$, a compact subset of \mathbb{R}^d . The function $m : \mathbb{R}^d \rightarrow \mathbb{R}$ is unknown, maybe complex, and it is assumed to be square-integrable.

Let \mathcal{D} be the set of densities,

$$\mathcal{D} = \left\{ \pi_\alpha : \pi_\alpha(x) = a_\alpha \exp(-|x|^\alpha), \text{ with } (a_\alpha)^{-1} = \int_{\mathbb{R}} \exp(-|x|^\alpha) dx, \alpha > 2 \right\}. \quad (2.16)$$

In this Chapter, we assume that the error term ε is equal to Z/σ_α , where Z is a random variable with density $\pi_\alpha \in \mathcal{D}$ and $\sigma_\alpha^2 = \text{var}(Z)$.

2.2.3 Main results

Let us begin with some notations, properties and assumptions that are needed to state [Result 1](#) and [Result 2](#).

Notations

- ✓ For a function $f \in \mathcal{H}$, let S_f be its support,

$$S_f = \{v \in \mathcal{P} : f_v \neq 0\}.$$

- ✓ Each kernel k_v , $v \in \mathcal{P}$ is associated with an integral operator T_{k_v} from $L^2(\mathcal{X}_v, P_v)$ to $L^2(\mathcal{X}_v, P_v)$ defined by:

$$\forall f \in L^2(\mathcal{X}_v, P_v), T_{k_v}(f) = \int_{\mathcal{X}_v} k_v(., t) f(t) dP_v(t).$$

For each $v \in \mathcal{P}$, let $\omega_{v,1} \geq \omega_{v,2} \geq \dots \geq 0$ be the eigenvalues of the integral operator T_{k_v} . Let us define the function $Q_{n,v}(t)$ for some positive t as follows:

$$Q_{n,v}(t) = \sqrt{\frac{5}{n} \sum_{\ell \geq 1} \min(t^2, \omega_{v,\ell})}.$$

- ✓ For some $\Delta > 0$ let $\nu_{n,v}$ be defined by:

$$\nu_{n,v} = \inf_t \left\{ Q_{n,v}(t) \leq \Delta t^2 \right\}. \quad (2.17)$$

For each $v \in \mathcal{P}$, $\nu_{n,v}$ refers to the minimax optimal rate for $L^2(\mathcal{X}, P_X)$ -estimation in the RKHS \mathcal{H}_v (Mendelson (2002)).

Remark 2.2.1 *The rate $\nu_{n,v}$, $v \in \mathcal{P}$, depends on the regularity of the RKHS via the decreasing rate of the eigenvalues $\{\omega_{v,\ell}\}_{\ell=1}^{\infty}$. When RKHS is of high regularity, i.e. when the eigenvalues $\{\omega_{v,\ell}\}_{\ell=1}^{\infty}$ decrease quickly, then the rate $\nu_{n,v}$, $v \in \mathcal{P}$ will be close to the parametric rate of convergence (see Section 3.3.1 of Chapter 3).*

Properties

The RKHS construction as described in Section 2.1.4.2 insures that the following properties are satisfied:

- P1 For all $v \in \mathcal{P}$, the functions $f_v \in \mathcal{H}_v$ are centered and are square-integrable,

$$E_X(f_v(X_v)) = 0 \text{ and } E_X(f_v^2(X_v)) < \infty.$$

- P2 For all $v, v' \in \mathcal{P}$, $v \neq v'$, the functions $f_v \in \mathcal{H}_v$ and $f_{v'} \in \mathcal{H}_{v'}$ are orthogonal with respect to $L^2(\mathcal{X}, P_X)$,

$$E_X(f_v(X_v) f_{v'}(X_{v'})) = 0.$$

Assumptions

- A1 For all $v \in \mathcal{P}$, the functions $f_v \in \mathcal{H}_v$ are uniformly bounded,

$$\exists R > 0 \text{ such that } \|f_v\|_{\infty} = \sup_{X_v} |f_v(X_v)| \leq R.$$

This assumption is satisfied as soon as the kernel k_v is bounded on the compact set \mathcal{X} . Indeed,

$$\|f_v\|_{\infty} \leq \sup_{X \in \mathcal{X}} \sqrt{k_v(X_v, X_v)} \|f_v\|_{\mathcal{H}_v}.$$

For each $v \in \mathcal{P}$, let us consider the quantity $\lambda_{n,v}$ defined by:

$$\lambda_{n,v} = \max\left(\nu_{n,v}, \sqrt{\frac{d}{n}}\right). \quad (2.18)$$

The regularization parameters μ_v and γ_v involved in the criterion (2.10) are chosen as follows:

A2 For some constant $C_1 > 10 + 4\Delta$,

$$\forall v \in \mathcal{P}, \mu_v = C_1 \lambda_{n,v}^2, \gamma_v = C_1 \lambda_{n,v}.$$

A3 There exists positive constants C_2, C_3 , and $0 < \beta < 1/\alpha$ such that the following conditions are satisfied:

$$\forall v \in \mathcal{P}, n\lambda_{n,v}^2 \geq -C_2 \log \lambda_{n,v}, \quad (2.19)$$

and

$$\forall f \in \mathcal{F}, \sum_{v \in S_f} \lambda_{n,v}^2 \leq C_3 n^{2\beta-1}. \quad (2.20)$$

Chapter 3 presents the two following results.

Result 1: upper bound of the empirical L^2 risk of the RKHS ridge group sparse estimator

Consider the regression model described in Section 2.2.2 with $\sigma = 1$. Let $\{(Y_i, X_i)\}_{i=1}^n$ be a n -sample with the same law as (Y, X) , and let $\{\varepsilon_i\}_{i=1}^n$ be the random errors that are independent and identically distributed (i.i.d.) like ε . Let also the RKHS ridge group sparse estimator \hat{f} be defined by (2.12) with $r_v = 1$ in (2.11). Under the assumptions A1, A2, and A3, there exists a positive constant C and $0 < \eta < 1$ (η tends to 0 as n increases) such that,

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\},$$

with probability greater than $1 - \eta$.

Let us comment on this [Result 1](#):

- R1 *Let f' be the function in \mathcal{F} such that the infimum of the right hand side of the oracle inequality is realized. The term $\|m - f'\|_n^2$ is the usual bias term. It quantifies both the approximation properties of the RKHS \mathcal{H} , and the bias-variance trade-off.*
- R2 *This result is similar to the one obtained in the Gaussian regression model at the cost of the additional Assumption (2.20). This assumption allows to obtain the same rate of convergence for the RKHS ridge group sparse estimator as in the Gaussian regression model (see [Huet and Taupin \(2017\)](#)). However,*

it implies some restrictions on the regularity of the RKHS \mathcal{H} . Indeed, as for all $v \in \mathcal{P}$, $\lambda_{n,v} \geq \nu_{n,v}$ (see Equation (2.18)), it follows that $\sum_{v \in S_f} \nu_{n,v}^2 \leq C_3 n^{2\beta-1}$, which implies some restrictions on the regularity of the RKHS: if β is small, which will be the case if α is large, then the RKHS should be of high regularity.

R3 By Equation (2.18), we also have that for all $v \in \mathcal{P}$, $\lambda_{n,v} \geq \sqrt{d/n}$. This assumption allows to control the probability of the $|\mathcal{P}|$ events (see Equation (3.48) of Chapter 3), where $\log(|\mathcal{P}|)$ is of order d .

R4 The Result 1 can be generalized to the case where $\sigma \neq 1$ in Equation (2.1), and where $r_v \neq 1$ in Equation (2.11). We refer to Remark 3.3.5 of Chapter 3 for a brief demonstration of this point.

Result 2: upper bound of the L^2 risk of the RKHS ridge group sparse estimator

Under the same assumptions as Result 1, we have with high probability for some positive constant C' that,

$$\|m - \hat{f}\|_2^2 \leq C' \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \|m - f\|_2^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\},$$

Remark 2.2.2 The Result 2 can be generalized to the case where $\sigma \neq 1$ in Equation (2.1), and where $r_v \neq 1$ in Equation (2.11) (see Remark 3.3.6 of Chapter 3 for more details about this point).

Rate of convergence

Under the same assumptions as Result 1, we have

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \left(\sum_{v \in S_f} \nu_{n,v}^2 + \frac{d|S_f|}{n} \right) \right\}.$$

This inequality highlights that the upper bound is relevant when the infimum is reached for functions f that have a sparse decomposition in \mathcal{H} , i.e. $|S_f|$ is small, and when d is small face to n . When d is large, the decomposition of functions in \mathcal{H} should be limited to interactions of a limited order, so that the number of elements in the estimated meta-model is of order smaller than d^r for some small r , say $r = 2$ for example. In such a case, the cardinality of \mathcal{P} will be smaller than d^2 . As we mentioned in Remark R3, the assumption $\lambda_{n,v} \geq \sqrt{d/n}$ is needed to control the value $\log(|\mathcal{P}|)$, which will be now smaller than $2 \log(d)$. Therefore, the value d in the definition of $\lambda_{n,v}$ (see Equation (2.18)) as well as the term $d|S_f|/n$ in the infimum above will be replaced by $2 \log(d)$ and $2 \log(d)|S_f|/n$, respectively.

2.2.4 Related works

Several authors studied the theoretical properties of estimators similar to the RKHS ridge group sparse estimator. Let us briefly review their framework and their results.

Meier et al. (2009) considered an estimator similar to the RKHS ridge group sparse estimator. Instead of adding two separate sparsity and smoothness penalties, they combine the two terms into a single sparsity and smoothness penalty. In the fixed design regression framework with error ε that is distributed as a sub-Gaussian random variable, they established the empirical risk upper bounds for the estimation of m onto the set of univariate additive functions. Afterwards, Raskutti et al. (2012) showed (in Section 3.4. of their paper) that the convergence rate of this estimator is sub-optimal.

Koltchinskii and Yuan (2010) considered a ridge group sparse type estimator defined on a set of additive functions where each term belongs to a RKHS. They do not assume that the input variables X_1, \dots, X_d are independent, nor that there is orthogonality between the RKHS spaces. Instead, they introduce some characteristics related to the degree of *dependence* of the RKHS spaces which insures *almost* orthogonality between these spaces. Under a global boundedness condition, they established upper bounds on the excess risk of their estimator by assuming that the function m has a sparse representation. A global boundedness condition means that the quantity $\sup_{f \in \mathcal{H}} \sup_{X \in \mathcal{X}} |f(X)|$ is assumed to be bounded independently of dimension d . Their results are valid for a large class of loss functions called *losses of quadratic type* which satisfy some defined boundedness conditions on the support of the output variable Y . Section 2.1. of their paper provides several examples of the framework for applying their results. Note that, the quadratic loss function in the case where Y is non-bounded does not belong to the class of the *losses of quadratic type*. The proofs of their results rely on the elementary empirical and Rademacher process methods such as symmetrization and concentration inequalities for Rademacher processes and Bernstein type exponential bounds.

Raskutti et al. (2012) assumed that the function m has a sparse univariate additive representation (as defined in Equation (2.8)) such that each univariate function lies in a RKHS. They proposed the ridge group sparse procedure to calculate the estimator of m , and studied the theoretical properties of their estimator in the Gaussian regression framework. They provided upper bounds for the empirical L^2 and the L^2 risks and a lower bound for the L^2 risk of their estimator over spaces of sparse additive models, including polynomials, splines and Sobolev classes.

Huet and Taupin (2017) studied the theoretical properties of the RKHS ridge group sparse estimator, in the Gaussian regression framework. They derived upper bounds of the empirical L^2 risk and the L^2 risk of the RKHS ridge group sparse estimator, i.e. the upper bounds with respect to the L^2 -norm and the empirical L^2 -norm for the distance between the true function m and its estimation \hat{f} into the RKHS \mathcal{H} .

Raskutti et al. (2012) and Huet and Taupin (2017) do not assume the global boundedness condition. Instead, they consider Assumption A1 where for all $v \in \mathcal{P}$,

$\sup_{X_v} |f_v(X_v)|$ is bounded. The proofs of their results rely on the probabilistic methods of the empirical Gaussian processes such as concentration inequalities and Sudakov minoration (Pisier (1989), Massart (2000), van de Geer et al. (2000), Ledoux (2001)), as well as the results on the Rademacher complexity of kernel classes (Mendelson (2002), Bartlett et al. (2005)).

2.2.5 Technical tools for the proofs

In this work, the upper bounds of the empirical L^2 risk and the L^2 risk of the RKHS ridge group sparse estimator are provided, in the regression framework where the error ε is non-Gaussian and non-bounded, and by considering a penalized least-squares criterion. In this case the conditions assumed in Koltchinskii and Yuan (2010) are not satisfied, and the usual probabilistic methods of the empirical Gaussian processes such as concentration inequalities and Sudakov minoration do not apply. The proofs of our results require different mathematical tools from those used in the past works:

- ✓ a Sudakov type minoration that for the non-Gaussian and non-bounded random variables,
- ✓ a concentration bound for the lower and upper tails of a convex function of the non-Gaussian and non-bounded random variables.

To the best of our knowledge, in our context of non-Gaussian and non-bounded errors, and with the least-squares criterion, the only Sudakov type minoration which allows to obtain the same rate of convergence for the RKHS ridge group sparse estimator as in the Gaussian regression framework (see Huet and Taupin (2017)), is the one obtained by Talagrand (1994). The minoration obtained by Talagrand (1994) is specific to the densities π_α (see Equation (2.16)). This is the reason why the class of densities \mathcal{D} is considered in this work. The Sudakov type minoration adapted to our work is derived from Theorem 3.1. in Talagrand (1994). Let us recall this minoration.

Let $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ be i.i.d. random variables distributed with density $\pi_\alpha \in \mathcal{D}$, and for a function $g : \mathbb{R}^{|v|} \mapsto \mathbb{R}$, $v \in \mathcal{P}$ belonging to a class of functions \mathcal{G} , let $V_{n,\varepsilon}$ be the empirical process associated with the random vector ε ,

$$V_{n,\varepsilon}(g) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i g(X_{v,i}). \quad (2.21)$$

Then, for all $\delta > 0$,

$$\begin{aligned} \frac{1}{K} \log N(\delta, \mathcal{G}, \|\cdot\|) &\leq \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^2 \mathbf{1}_{[2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|, \infty)}(\delta) \\ &+ \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^\alpha \mathbf{1}_{(0, 2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|]}(\delta), \end{aligned} \quad (2.22)$$

where K is a constant that depends on α only, $N(\delta, \mathcal{G}, \|\cdot\|)$ is the δ -covering number of the metric space $(\mathcal{G}, \|\cdot\|)$, and $1_A : \mathcal{A} \rightarrow \{0, 1\}$ is the indicator function of $A \subset \mathcal{A}$,

$$1_A(a) = \begin{cases} 1 & \text{if } a \in A, \\ 0 & \text{if } a \notin A. \end{cases}$$

Concerning the concentration bound, it is shown that the distribution functions associated with the densities $\pi_\alpha \in \mathcal{D}$ belong to a class of distribution functions defined by Adamczak (2005), for which the log-Sobolev inequality (Gross (1975)) is satisfied (see Lemma 3.4.2 in Chapter 3). Shu and Strzelecki (2017) provided bounds for the lower and upper tails of convex functions of independent random variables which satisfy the log-Sobolev inequality. As the distribution functions associated with the densities $\pi_\alpha \in \mathcal{D}$ satisfy the log-Sobolev inequality, the concentration inequality derived by Shu and Strzelecki (2017) holds for them. The concentration inequality adapted to our work is derived from Corollary 1.7. in Shu and Strzelecki (2017) (see Corollary 3.4.2 in Chapter 3).

2.3 Summary of Chapter 4

2.3.1 Objectives and results

An R package, called **RKHSMetaMod**, is developed to implement the ridge group sparse procedure described in Section 2.1.5.1. This package allows to:

- ✓ calculate reproducing kernels as described in Section 2.1.4.2, and their associated Gram matrices,
- ✓ implement the RKHS ridge group sparse procedure and a special case of it called RKHS group lasso procedure (when $\gamma_v = 0$, $v \in \mathcal{P}$ in criterion (2.13)) in order to estimate the terms f_v^* in the Hoeffding decomposition of f^* leading to an estimation of the function m ,
- ✓ choose the tuning parameters $\mu_v, \gamma_v, v \in \mathcal{P}$ in the criterion (2.13) using a procedure that leads to obtain the *best* RKHS ridge group sparse estimator in terms of the prediction quality,
- ✓ estimate the Sobol indices of the function m as described in Section 2.1.5.2.

The **RKHSMetaMod** package provides an interface from R statistical computing environment to the C++ libraries **Eigen** and **GSL**. In order to optimize the execution time and the storage memory, except for a function that is written in R, all of the functions of this package are written using the efficient C++ libraries through **RcppEigen** and **RcppGSL** packages. These functions are then interfaced in the R environment in order to propose an user friendly package. The **RKHSMetaMod** package is dedicated to the meta-model estimation on a RKHS. The convex optimization algorithms used in this package are adapted to take into account the problem of high dimensionality in this context. This package is available from

the Comprehensive R Archive Network (CRAN) at <https://cran.r-project.org/web/packages/RKHSMetaMod/>.

2.3.2 Presentation of the model

Consider a phenomenon described by a model m depending on d input variables $X = (X_1, \dots, X_d)$. This model m from \mathbb{R}^d to \mathbb{R} may be a known model that is calculable in all points X , or a regression model as defined in Equation (2.1). In the second case, the error ε is assumed to be centered with a finite variance. The components of X are independent and uniformly distributed on $\mathcal{X} = [0, 1]^d$, i.e. $X \sim P_X = P_1 \times \dots \times P_d$, with P_a , $a = 1, \dots, d$ being the uniform law on the interval $[0, 1]$. The model m may present high complexity as strong non-linearities and high order interaction effects, and it is assumed to be square-integrable.

2.3.3 Criterion to minimize

Let us consider the parametric form of the RKHS ridge group sparse criterion defined in Equation (2.13), where γ_v and μ_v , $v \in \mathcal{P}$ are chosen as follows:

For each $v \in \mathcal{P}$, let γ'_v and μ'_v be the weights that are chosen suitably. Then,

$$\gamma_v = \gamma \times \gamma'_v \text{ and } \mu_v = \mu \times \mu'_v \text{ with } \gamma, \mu \in \mathbb{R}^+.$$

Remark 2.3.1 *This formulation simplify the choice of the tuning parameters, since instead of tuning the parameters γ_v and μ_v for all $v \in \mathcal{P}$, only two parameters γ and μ are tuned. Moreover, the weights γ'_v and μ'_v , $v \in \mathcal{P}$, may be of interest in applications. For example, one can take weights that increase with the cardinal of v in order to favour effects with small interaction order between variables.*

For the sake of simplicity, in the rest of this Chapter for all $v \in \mathcal{P}$ the weights γ'_v and μ'_v are assumed to be setted as 1, and the RKHS ridge group sparse criterion is then expressed as follows:

$$C(f_0, \theta) = \left\{ \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \gamma \sum_{v \in \mathcal{P}} \|K_v \theta_v\| + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\| \right\}.$$

By considering only the second part of the penalty function in the criterion above, i.e. set $\gamma = 0$, we obtain the RKHS group lasso criterion,

$$C_g(f_0, \theta) = \left\{ \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\| \right\},$$

which is a group lasso criterion (Yuan and Lin (2006)) up to a scale transformation.

We would like to mention that, in the **RKHSMetaMod** package, the solutions of the RKHS group lasso algorithm is used in order to initialize the input parameters of the RKHS ridge group sparse algorithm. Indeed, the penalty function in the RKHS

group lasso criterion $C_g(f_0, \theta)$ insures the sparsity in the solution. Therefore, for a given value of μ , by implementing the RKHS group lasso algorithm, we obtain a solution with few terms in its additive decomposition.

From now on, the tuning parameter in the RKHS group lasso algorithm will be denoted by:

$$\mu_g = \sqrt{n}\mu. \quad (2.23)$$

2.3.3.1 Choice of the tuning parameters

While dealing with an optimization problem, one of the essential steps is to choose appropriately the tuning parameters. To do so,

- ✓ first, a grid of values of the tuning parameters μ and γ is chosen.

Let μ_{\max} be the smallest value of μ_g (see Equation (2.23)), such that the solution to the minimization of the RKHS group lasso problem for all $v \in \mathcal{P}$ is $\theta_v = 0$. We have,

$$\mu_{\max} = \max_v \left(\frac{2}{\sqrt{n}} \|K_v^{1/2}(Y - \bar{Y})\| \right).$$

In order to set up the grid of values of μ , one may find μ_{\max} , and then a grid of values of μ is defined as follows:

$$\mu_l = \frac{\mu_{\max}}{(\sqrt{n} \times 2^l)}, \quad l \in \{1, \dots, l_{\max}\}.$$

The grid of values of γ is chosen by the user.

- ✓ next, for the grid of values of μ and γ a sequence of estimators is calculated. Each estimator associated with the pair (μ, γ) in the grid of values of μ and γ , denoted by $\hat{f}_{(\mu, \gamma)}$, is the solution of the RKHS ridge group sparse optimization problem or the RKHS group lasso optimization problem if $\gamma = 0$.
- ✓ finally, the obtained estimators $\hat{f}_{(\mu, \gamma)}$ are evaluated using a testing dataset,

$$\{(Y_i^{\text{test}}, X_i^{\text{test}})\}_{i=1}^{n^{\text{test}}}.$$

The prediction error associated with the estimator $\hat{f}_{(\mu, \gamma)}$ is calculated by,

$$\text{ErrPred}(\mu, \gamma) = \frac{1}{n^{\text{test}}} \sum_{i=1}^{n^{\text{test}}} (Y_i^{\text{test}} - \hat{f}_{(\mu, \gamma)}(X_i^{\text{test}}))^2,$$

where for $S_{\hat{f}}$ being the support of the estimator $\hat{f}_{(\mu, \gamma)}$,

$$\hat{f}_{(\mu, \gamma)}(X^{\text{test}}) = \hat{f}_0 + \sum_{v \in S_{\hat{f}}} \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v^{\text{test}}).$$

The pair $(\hat{\mu}, \hat{\gamma})$ with the smallest value of the prediction error is chosen, and the estimator $\hat{f}_{(\hat{\mu}, \hat{\gamma})}$ is considered as the *best* estimator of the function m , in terms of the prediction error.

In the **RKHSMetaMod** package, the algorithm to calculate a sequence of the estimators \hat{f} , the value of μ_{\max} , and the prediction error are implemented as `RKHSMetMod`, `mu_max`, and `PredErr` functions, respectively.

2.3.3.2 Estimation of the Sobol indices

The Sobol indices of the function m are estimated by the empirical Sobol indices of the estimator \hat{f} as described in Section 2.1.5.2,

$$\hat{S}_v = \begin{cases} \frac{\widehat{\text{var}}(\hat{f}_v(X_v))}{\sum_{v \in \mathcal{P}} \widehat{\text{var}}(\hat{f}_v(X_v))} & \text{for } v \in S_{\hat{f}}, \\ 0 & \text{for } v \notin S_{\hat{f}}. \end{cases}$$

In the **RKHSMetaMod** package, the algorithm to calculate the empirical Sobol indices \hat{S}_v , $v \in \mathcal{P}$ is implemented as `SI_emp` function.

2.3.4 Algorithms

The **RKHSMetaMod** package implements two optimization algorithms: the RKHS ridge group sparse and the RKHS group lasso. These algorithms rely on the Gram matrices K_v , $v \in \mathcal{P}$, that have to be positive definite. Therefore, the first and essential step in the **RKHSMetaMod** package, is to calculate these matrices and insure their positive definiteness.

The second step is to calculate the estimator \hat{f} . In the **RKHSMetaMod** package two different objectives based on different procedures are considered in order to calculate this estimator:

- ✓ The estimator with the *best* prediction quality:

In this case the *best* estimator is calculated using the procedure as described in Section 2.3.3.1.

- ✓ The estimator with at most *qmax* active groups:

The tuning parameter γ is set as zero. A value of μ for which the number of groups in the solution of the RKHS group lasso problem is equal to *qmax*, is computed. This value will be denoted by μ_{qmax} . Then, the RKHS ridge group sparse algorithm is implemented for a grid of values of $\gamma \neq 0$ and the value μ_{qmax} .

This procedure is implemented in the **RKHSMetaMod** package as `RKHSMetMod_qmax` function.

2.3.4.1 Calculation of the Gram matrices

The available kernels in the **RKHSMetaMod** package are: linear kernel, quadratic kernel, brownian kernel, matern kernel and gaussian kernel. The choice of the kernel that is done by the user, determines the functional approximation space. For a chosen kernel, the algorithm to calculate the Gram matrices K_v , $v \in \mathcal{P}$ in

the **RKHSMetaMod** package is implemented as `calc_Kv` function, and is based on three essential points:

- ✓ Modify the chosen kernel:

In order to satisfy the conditions of constructing the RKHS \mathcal{H} described in Section 2.1.4.2, these kernels are modified according to Equation (2.7). Let us take the example of the Brownian kernel:

Example 2.3.1 *The usual presentation of the brownian kernel is as follows:*

$$k_a(X_a, X'_a) = \min(X_a, X'_a) + 1.$$

The RKHS associated with the kernel k_a is the set,

$$\mathcal{H}_a = \left\{ f : [0, 1] \rightarrow \mathbb{R} \text{ is absolutely continuous, and } f(0) = 0, \int_0^1 f'(X_a)^2 dX_a < \infty \right\},$$

with the inner product

$$\langle f, h \rangle_{\mathcal{H}_a} = \int_0^1 f'(X_a)h'(X_a)dX_a.$$

The kernel k_{0a} associated with the brownian kernel is calculated as follows,

$$\begin{aligned} k_{0a} &= \min(X_a, X'_a) + 1 - \frac{(\int_0^1 (\min(X_a, U) + 1)dU)(\int_0^1 (\min(X'_a, U) + 1)dU)}{(\int_0^1 \int_0^1 (\min(U, V) + 1)dUdV)}, \\ &= \min(X_a, X'_a) + 1 - \frac{3}{4} \left(1 + X_a - \frac{X_a^2}{2}\right) \left(1 + X'_a - \frac{X'^2_a}{2}\right). \end{aligned}$$

The RKHS associated with the kernel k_{0a} is the set,

$$\mathcal{H}_{0a} = \left\{ f \in \mathcal{H}_a : \int_0^1 f(X_a)dX_a = 0 \right\}.$$

Finally, the RKHS $\mathcal{H} = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v$ is the following set,

$$\mathcal{H} = \left\{ f : [0, 1]^d \rightarrow \mathbb{R} : f = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \text{ with } f_v \in \mathcal{H}_v \right\}.$$

- ✓ Calculate the Gram matrices K_v for all v :

First, for all $a = 1, \dots, d$ the Gram matrices K_a associated with kernels k_{0a} are calculated using Equation (2.7),

$$(K_a)_{i,i'} = k_{0a}(X_{ai}, X_{ai'}).$$

Then, for all $v \in \mathcal{P}$, the Gram matrices K_v associated with kernel $k_v = \prod_{a \in v} k_{0a}$ are calculated as follows:

$$K_v = \bigodot_{a \in v} K_a,$$

where \bigodot denotes the Hadamard product.

- ✓ Insure the positive definiteness of the matrices K_v , $v \in \mathcal{P}$:

The output of the function `calc_Kv` is one of the input arguments of the functions associated with the RKHS group lasso and the RKHS ridge group sparse algorithms. As both of these algorithms rely on the positive definiteness of these matrices, it is mandatory to have K_v , $v \in \mathcal{P}$ that are positive definite. For this reason, the `calc_Kv` function modifies the eigenvalues of the matrices K_v , $v \in \mathcal{P}$ if necessary.

For each group $v \in \mathcal{P}$, let $\lambda_{v,\max}$ and $\lambda_{v,\min}$ be respectively the maximum and the minimum eigenvalues associated with the matrix K_v , and let "tol" be a positive scalar to be fixed. For each matrix K_v ,

$$\text{"if } \lambda_{v,\min} < \lambda_{v,\max} \times \text{tol"},$$

then, $\lambda_{v,\max} \times \text{tol}$ is added to all eigenvalues of K_v .

The value of "tol" is set as $1e^{-8}$ by default, but one may consider a smaller or greater value for it depending on the kernel chosen and the value of n .

2.3.4.2 Optimization algorithms

RKHS group lasso In order to solve the RKHS group lasso optimization problem, the classical block coordinate descent algorithm is used (Boyd et al. (2011), Bubeck (2015)). The minimization of criterion $C_g(f_0, \theta)$ is done along each group v at a time. At each step of the algorithm, the criterion is minimized as a function of the current block's parameters, while the parameters values for the other blocks are fixed to their current values. The procedure is repeated until convergence.

In the **RKHSMetaMod** package the classical block coordinate descent algorithm to solve the RKHS group lasso optimization problem is implemented as `RKHSgrplasso` function.

RKHS ridge group sparse In order to solve the RKHS ridge group sparse optimization problem, an adapted block coordinate descent algorithm is proposed. This algorithm provides two steps:

- Step 1** Initialize the input parameters by the solutions of the RKHS group lasso algorithm for each value of the tuning parameter μ , and run the RKHS ridge group sparse algorithm through active support of the RKHS group lasso solutions until it achieves convergence.

This step is provided in order to decrease the execution time.

- Step 2** Re-initialize the input parameters with the obtained solutions of **Step 1** and implement the RKHS ridge group sparse algorithm through all groups in \mathcal{P} until it achieves convergence.

This second step makes it possible to verify that no group is missing in the output of **Step 1**.

The adapted block coordinate descent algorithm to solve RKHS ridge group sparse optimization problem is implemented in the **RKHSMetaMod** package, as `pen_MetMod` function.

2.4 Summary and perspectives

The work presented in this thesis focuses on the problem of estimating a meta-model that approximates the Hoeffding decomposition of a complex model, denoted m . The model m depends on d input variables X_1, \dots, X_d that are independent and have a known law. The meta-model belongs to a RKHS \mathcal{H} , which is constructed in a way such that the additive decomposition of any function f in \mathcal{H} is the Hoeffding decomposition of f (Durrande et al. (2013)). The estimator of the meta-model, denoted \hat{f} , minimizes a least-squares criterion penalized by a penalty function which is the sum of the Hilbert norm and the empirical L^2 -norm. This procedure, called RKHS ridge group sparse, allows both to select and estimate the terms in the Hoeffding decomposition of the meta-model, and therefore, to select the Sobol indices that are non-zero and estimate them (Huet and Taupin (2017)).

The first part of this work is dedicated to study the theoretical properties of the estimator \hat{f} in the regression framework where the error ε that is non-Gaussian and non-bounded. The upper bounds of the empirical L^2 and the L^2 risks of this estimator are provided.

In the second part of this work, the procedure of calculating \hat{f} is implemented in an R package, called **RKHSMetaMod**. In order to optimize the execution time and also the storage memory, except for a function that is written in R, all of the functions in this package are written using C++ libraries **GSL** and **Eigen**. They are then interfaced with the R environment in order to propose an user friendly package. The **RKHSMetaMod** package deals both with a calculable model and a regression model. A simulation study is provided in order to validate the performance of the package functions in terms of the predictive quality of the estimator obtained and the estimation of the Sobol indices.

Like all research works that are carried out in a limited period of time, many pistes have not been explored in this work and there are several perspectives to be considered for further study. Let us mention some of them.

2.4.1 Non-independent input variables

In both parts of this thesis, the input variables X_1, \dots, X_d are assumed to be independent and their law is known. Under these assumptions, it is possible to construct the approximation spaces such that any function in these spaces is decomposed according to its Hoeffding decomposition. This decomposition is unique and the terms of it are orthogonal.

If the variables X_1, \dots, X_d are not independent, there is no longer orthogonality between the terms of the decomposition on the approximation spaces and the decomposition of a function on these spaces is not necessarily unique. It follows that

the decomposition of the variance given in Equation (2.4) is no longer valid, nor the calculation of the Sobol indices. However, approximating a model on a functional space by an additive decomposition may be interesting in practice, since the estimation of the meta-model can still help the interpretation of the effects of the input variables on the output variable.

The case where the variables X_1, \dots, X_d are not independent has been considered by Koltchinskii and Yuan (2010). Their approximation space \mathcal{H} is the linear span (l.s.) of a large dictionary consisting of N RKHS spaces $\mathcal{H}_1, \dots, \mathcal{H}_N$,

$$\mathcal{H} = \text{l.s.} \bigcup_{j=1}^N \mathcal{H}_j.$$

The space \mathcal{H} consists of all functions f that have an additive representation of the form,

$$f = \sum_{j=1}^N f_j(X), \quad f_j \in \mathcal{H}_j, \quad j = 1, \dots, N. \quad (2.24)$$

Under some assumptions on the degree of *dependence* of the RKHS spaces \mathcal{H}_j spaces ensuring *almost* orthogonality between these spaces, Koltchinskii and Yuan (2010) established upper bounds in the excess risk of a ridge group sparse type estimator.

Note that, the approximation spaces considered in this thesis are a special case of the spaces considered by Koltchinskii and Yuan (2010). However, the context in which their results are established differs from ours in several points. On the one hand, the functions in the space \mathcal{H} are assumed to be uniformly bounded. On the other hand, the statistical model is different and in particular the case of the regression model with non-bounded additive error is not considered in their work.

My objective would be then to establish the risk upper bounds of a ridge group sparse estimator on the approximation spaces constructed as proposed by Durrande et al. (2013), in a context where the input variables X_1, \dots, X_d are non-independent, for the regression model with non-bounded additive error. One of the essential steps of the proof relies on the upper bounding the L^2 -norm in \mathcal{H} by the empirical L^2 -norm in \mathcal{H} (Lemma 3.5.4). The calculation of this upper bound requires the control of the moments of order 4 of the functions in \mathcal{H} . In the case where the variables X_1, \dots, X_d are independent, this control is obtained since there is orthogonality between the terms in the decomposition of the functions in \mathcal{H} . In the contrary case, the assumptions on the degree of *dependence* of \mathcal{H}_j spaces formulated by Koltchinskii and Yuan (2010) is not sufficient to handle the calculation of moments of order 4. It remains therefore to continue this work to establish the upper bounds of the risk of the RKHS ridge group sparse estimator under assumptions which have to be specified. To our knowledge, this case has not been studied until now.

Concerning the implementation of sensitivity analysis in this case, as the calculation of the Sobol indices is not possible any more, one may consider Shapley values (Shapley (1953)), see for example Owen (2014), Song et al. (2016), Owen and Prieur

(2017), Benoumechiara and Elie-Dit-Cosaque (2019), Broto et al. (2019), Iooss and Prieur (2019).

2.4.2 Generalization to the regression framework with log-concave error

Result 1 shows that the risk upper bound in the regression setting with errors that are distributed with density $\pi_\alpha \in \mathcal{D}$ is the same as the one obtained in the regression setting with Gaussian errors. However, the class of densities π_α is restrictive, and it would be interesting to obtain a result for larger density classes, such as log-concave densities for example.

As explained in Section 2.2.5, one of the essential steps of the proof of **Result 1** is based on a Sudakov type minoration of the expectation of the empirical process associated with the random variables that are assumed to be non-Gaussian and non-bounded.

In the Gaussian case, Sudakov's minoration is stated as follows (**Pisier (1989)**):

Let $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ be i.i.d. Gaussian random variables, and for a function $g : \mathbb{R}^{|v|} \mapsto \mathbb{R}$, $v \in \mathcal{P}$ belonging to a class of functions \mathcal{G} , let $V_{n,\varepsilon}$ be the empirical process associated with the random vector ε defined in Equation (2.21). Then, for all $\delta > 0$,

$$\frac{1}{C} \log N(\delta, \mathcal{G}, \|\cdot\|) \leq \left(\frac{n E_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^2, \quad (2.25)$$

where C is a constant, and $N(\delta, \mathcal{G}, \|\cdot\|)$ is the δ -covering number of the metric space $(\mathcal{G}, \|\cdot\|)$.

It remains then to characterize the complexity of the functional space \mathcal{G} to obtain a minoration of the expectation of the empirical process and to deduce the result in the risk bound.

When $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ are distributed according to density π_α , Theorem 3.1. in **Talagrand (1994)** establishes the inequality given in Equation (3.34) from which we could deduce the minoration of the expectation of the empirical process given in Equation (2.22).

When ε is non-Gaussian and non-bounded, a Sudakov type minoration for the independent log-concave random variables is given by **Latała (2014)**. A measure on \mathbb{R}^n with the full dimensional support is log-concave if and only if it has a density of the form $\exp(-\phi(x))$, where $\phi : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is convex (**Borell (1974)**). Let us recall the Sudakov type minoration obtained by **Latała (2014)**:

Let $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ be i.i.d. log-concave random variables, then:

$$\frac{1}{K'} \min \left(c^2 \delta, \log N(2 \times \max(c\delta^{1/2}, c^2 \delta), \mathcal{G}, \|\cdot\|) \right) \leq n E_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|, \quad (2.26)$$

where K' is a universal constant, and $c = 1/\max(512c', 8)$ for c' being a universal constant.

We could not deduce from inequality (2.26) the adapted Sudakov type minoration that leads to obtain the *optimal* rate of convergence for the RKHS ridge group sparse estimator. By *optimal* we mean the same rate of convergence as in the Gaussian regression setting (see [Huet and Taupin \(2017\)](#)). This is the reason why, in this work, the densities $\pi_\alpha \in \mathcal{D}$ are considered. Nevertheless, some additional work in that direction together with bibliography research is a worthwhile direction.

Risk upper bounds for RKHS ridge group sparse estimator in the regression model with non-Gaussian and non-bounded error

Abstract

We consider the problem of estimating a meta-model of an unknown regression model with non-Gaussian and non-bounded error. The meta-model belongs to a reproducing kernel Hilbert space constructed as a direct sum of Hilbert spaces leading to an additive decomposition including the variables and interactions between them. The estimator of this meta-model is calculated by minimizing an empirical least-squares criterion penalized by the sum of the Hilbert norm and the empirical L^2 -norm. In this context, the upper bounds of the empirical L^2 risk and the L^2 risk of the estimator are established.

Keywords: meta-model, reproducing kernel Hilbert space, ridge group sparse, risk upper bound.

3.1 Introduction

Let us consider the following regression model:

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0, \quad (3.1)$$

where the variables $X = (X_1, \dots, X_d)$ are independent with a known law $P_X = \bigotimes_{a=1}^d P_a$ on $\mathcal{X} = \prod_{a=1}^d \mathcal{X}_a$, a compact subset of \mathbb{R}^d . The number d of components of X may be large. The model m from \mathbb{R}^d to \mathbb{R} maybe complex, presenting strong non-linearities, and it is assumed to be square-integrable, i.e. $m \in L^2(\mathcal{X}, P_X)$.

Let \mathcal{D} be the set of densities,

$$\mathcal{D} = \left\{ \pi_\alpha : \pi_\alpha(x) = a_\alpha \exp(-|x|^\alpha), \text{ with } (a_\alpha)^{-1} = \int_{\mathbb{R}} \exp(-|x|^\alpha) dx, \alpha > 2 \right\}. \quad (3.2)$$

In this Chapter, we assume that the error term ε is equal to Z/σ_α , where Z is a random variable with density $\pi_\alpha \in \mathcal{D}$ and σ_α^2 is its variance, i.e. $\text{var}(Z) = \sigma_\alpha^2$.

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator 62 in the regression model with non-Gaussian and non-bounded error

Based on n data points $\{(X_i, Y_i)\}_{i=1}^n$, a meta-model that approximates the Hoeffding decomposition of m is estimated. This meta-model belongs to a reproducing kernel Hilbert space (RKHS), which is constructed as a direct sum of Hilbert spaces (Durrande et al. (2013)). The estimation of the meta-model is carried out via a penalized least-squares minimization allowing to select the subsets of variables X that contribute to predict the output Y (Huet and Taupin (2017)).

Let us be more precise on the Hoeffding decomposition. Let \mathcal{P} be the set of all the subsets of $\{1, \dots, d\}$ with dimension 1 to d , and for all $v \in \mathcal{P}$ and $X \in \mathcal{X}$, let X_v be the vector with components X_a for all $a \in v$. Let also $|A|$ be the cardinality of a set A and for all $v \in \mathcal{P}$, let $m_v : \mathbb{R}^{|v|} \rightarrow \mathbb{R}$ be a function of X_v . Then, the Hoeffding decomposition of m is written as (Hoeffding (1948), Sobol (1993), van der Vaart (1998)),

$$m(X) = m_0 + \sum_{v \in \mathcal{P}} m_v(X_v), \quad (3.3)$$

where m_0 is a constant.

This decomposition (3.3) is unique (Sobol (1993)), all the functions m_v are centered, and they are orthogonal with respect to $L^2(\mathcal{X}, P_X)$.

The Hoeffding decomposition of m is approximated by the orthogonal projection of m on a RKHS \mathcal{H} which is constructed as a direct sum of Hilbert spaces (Durrande et al. (2013)).

Let $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ be the inner product in \mathcal{H} , and let k and k_v be the reproducing kernels associated with the RKHS \mathcal{H} and the RKHS \mathcal{H}_v , respectively. The properties of the RKHS \mathcal{H} insures that any function $f \in \mathcal{H}$, $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ can be written as the following decomposition:

$$f(X) = \langle f, k(X, \cdot) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \quad (3.4)$$

where f_0 is a constant, and $f_v : \mathbb{R}^{|v|} \rightarrow \mathbb{R}$ is defined by,

$$f_v(X) = \langle f, k_v(X, \cdot) \rangle_{\mathcal{H}}.$$

For all $v \in \mathcal{P}$, the functions $f_v(X_v)$ are centered and for all $v, v' \in \mathcal{P}$, $v \neq v'$, the functions $f_v(X_v)$ and $f_{v'}(X_{v'})$ are orthogonal with respect to $L^2(\mathcal{X}, P_X)$. Therefore, the decomposition of any function f presented in Equation (3.4) is unique and is its Hoeffding decomposition.

The meta-model f^* that approximates the Hoeffding decomposition of m is defined as follows:

$$f^* = \arg \min_{f \in \mathcal{H}} \|m - f\|_2^2 = \arg \min_{f \in \mathcal{H}} E_X (m(X) - f(X))^2.$$

Since the function f^* belongs to the RKHS \mathcal{H} , its decomposition on \mathcal{H} is its Hoeffding decomposition:

$$f^* = f_0^* + \sum_{v \in \mathcal{P}} f_v^*. \quad (3.5)$$

And for all $v \in \mathcal{P}$, the function f_v^* in Equation (3.5) approximates the function m_v in Equation (3.3).

Decomposition (3.5) contains $|\mathcal{P}|$ terms f_v^* to be estimated. The cardinality of \mathcal{P} is equal to $2^d - 1$ which may be huge since it raises very quickly by increasing d . In order to deal with this problem, one may estimate f^* by a sparse estimator $\hat{f} \in \mathcal{H}$. To this purpose, the estimation of f^* is done on the basis of n observations by minimizing an empirical least-squares criterion penalized by the sum of the Hilbert norm and the empirical norm. This procedure, called ridge group sparse, estimates the groups v that are suitable for predicting f^* , and the relationship between f_v^* and X_v for each group v (Huet and Taupin (2017)). The estimator so obtained is called the RKHS ridge group sparse estimator.

Several authors studied the theoretical properties of estimators similar to the RKHS ridge group sparse estimator. Let us briefly review their framework and their results.

Meier et al. (2009) considered an estimator similar to the RKHS ridge group sparse estimator. Instead of adding two separate sparsity and smoothness penalties, they combine these two terms into a single sparsity and smoothness penalty. In the fixed design regression model with error ε that is distributed as a sub-Gaussian random variable, they established upper bounds of the empirical risk for estimating the projection of m onto the set of univariate additive functions. Afterwards, Raskutti et al. (2012) showed (in Section 3.4. of their paper) that the convergence rate of this estimator is sub-optimal.

Koltchinskii and Yuan (2010) considered a more general RKHS including the functions that have an additive representation over kernel spaces and obtained an estimator based on a ridge group sparse type procedure. Under a global boundedness condition, they established upper bounds on the excess risk assuming that the function m has a sparse representation. A global boundedness condition means that the quantity $\sup_{f \in \mathcal{H}} \sup_{X \in \mathcal{X}} |f(X)|$ is assumed to be bounded independently of dimension d . Their results are valid for a large class of loss functions, and for distributions of the observations Y such that some defined boundedness conditions on the loss functions are satisfied (see Section 2.1. of their paper). In their framework, the input variables X are not assumed to be independent and there is no orthogonality assumption between the kernel spaces. Instead, the authors introduced some characteristics related to the degree of *dependence* of their kernel spaces which insures *almost* orthogonality between these spaces. Their method to derive their upper bounds relies on the elementary empirical and Rademacher process methods such as symmetrization and concentration inequalities for Rademacher processes and Bernstein type exponential bounds.

Raskutti et al. (2012) assumed that the function m has a sparse univariate additive representation, i.e. $m = \sum_{a \in \mathcal{S}} m_a(X_a)$ for $m_a(X_a)$ being univariate functions and $|\mathcal{S}| < d$, such that each univariate function m_a lies in a RKHS \mathcal{H}_a . They used the ridge group sparse procedure to calculate the estimator of m , and studied the theoretical properties of their estimator in the Gaussian regression model, i.e. ε in Equation (3.1) is distributed as a centered Gaussian random variable. They pro-

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator 64 in the regression model with non-Gaussian and non-bounded error

vided upper bounds for the integrated and the empirical risks and a lower bound for the integrated risk of their estimator over spaces of sparse additive models, including polynomials, splines and Sobolev classes.

Huet and Taupin (2017) studied the theoretical properties of the RKHS ridge group sparse estimator in the Gaussian regression model. They derived upper bounds with respect to the L^2 -norm and the empirical L^2 -norm for the distance between the true function m and its estimation \hat{f} into the RKHS \mathcal{H} .

Raskutti et al. (2012) and Huet and Taupin (2017) did not assume the global boundedness condition. Instead, they assumed that each function within the unit ball of the Hilbert space \mathcal{H}_v is uniformly bounded by a constant. The proof of their results is based on the probabilistic methods of empirical Gaussian process such as concentration inequalities and Sudakov minoration (e.g. Pisier (1989), Massart (2000), van de Geer et al. (2000), Ledoux (2001)), as well as results on the Rademacher complexity of kernel classes (Mendelson (2002), Bartlett et al. (2005)).

In this Chapter, the upper bounds of the empirical L^2 risk and the L^2 risk of the RKHS ridge group sparse estimator are provided, in the regression model (see Equation (3.1)) with non-Gaussian and non-bounded error ε , and by considering a quadratic loss function. In this case the conditions assumed in Koltchinskii and Yuan (2010) are not satisfied, and the empirical Gaussian process methods such as concentration inequalities and Sudakov minoration can not be used.

The proof of our results requires different mathematical tools than those used in the works mentioned above:

- ✓ a Sudakov type minoration that is satisfied for the non-Gaussian and non-bounded random variables,
- ✓ a concentration bound for the lower and upper tails of a convex function of the random variables $\{\varepsilon_i\}_{i=1}^n$ that are non-Gaussian and non-bounded.

To the best of our knowledge, in our context of regression model with non-Gaussian and non-bounded error ε , and with quadratic loss function, the only Sudakov type minoration which allows to obtain the same rate of convergence for the RKHS ridge group sparse estimator as in the Gaussian regression model (see Huet and Taupin (2017)), is the one obtained by Talagrand (1994). The minoration obtained by Talagrand (1994) is specific to the densities $\pi_\alpha \in \mathcal{D}$ as defined in Equation (3.2). This is the reason why this class of densities is considered in this work.

Concerning the concentration bound, it can be shown that the distribution functions associated with the densities $\pi_\alpha \in \mathcal{D}$ belong to a class of distribution functions defined by Adamczak (2005), for which the log-Sobolev inequality (Gross (1975)) is satisfied. Shu and Strzelecki (2017) provided bounds for the lower and upper tails of convex functions of independent random variables which satisfy the log-Sobolev inequality. Since the distribution functions associated with the densities $\pi_\alpha \in \mathcal{D}$ satisfy the log-Sobolev inequality, the concentration inequality derived by Shu and Strzelecki (2017) holds for them.

This Chapter is organised as follows: The RKHS construction and the procedure for estimating a meta-model are presented in Section 3.2. The theoretical properties of the RKHS ridge group sparse estimator are stated in Theorem 3.3.1 and Corollary 3.3.2. The proof of Theorem 3.3.1 is postponed in Section 3.5. In Section 3.4 the main arguments of the proof of Theorem 3.3.1 and motivation for the choice π_α are detailed.

3.2 Meta-modelling and the RKHS ridge group sparse estimator

The independency between the input variables X allows to write the function m according to its Hoeffding decomposition presented in Equation (3.3),

$$m(X) = m_0 + \sum_{v \in \mathcal{P}} m_v(X_v).$$

The unknown function m is approximated by its orthogonal projection, denoted f^* , on a RKHS, denoted \mathcal{H} , that is constructed as a direct sum of Hilbert spaces. The RKHS \mathcal{H} is associated with a so-called ANOVA kernel which is defined in order to obtain the analytical expression of the terms of the Hoeffding decomposition of the functions of \mathcal{H} . As f^* is the orthogonal projection of m on \mathcal{H} , each term in its decomposition is an approximation of the associated term in the Hoeffding decomposition of m . The construction of the RKHS \mathcal{H} has been proposed by [Durrande et al. \(2013\)](#) that we recall briefly in the following.

3.2.1 RKHS construction

Let $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ be a subset of \mathbb{R}^d . For each $a \in \{1, \dots, d\}$, we choose a RKHS \mathcal{H}_a , and its associated kernel k_a defined on the set $\mathcal{X}_a \subset \mathbb{R}$ such that the two following properties are satisfied:

- (i) $k_a : \mathcal{X}_a \times \mathcal{X}_a \rightarrow \mathbb{R}$ is $P_a \otimes P_a$ measurable,
- (ii) $E_{X_a} \sqrt{k_a(X_a, X_a)} < \infty$.

The property (ii) depends on the kernel k_a , $a = 1, \dots, d$ and the distribution of X_a , $a = 1, \dots, d$. It is not very restrictive since it is satisfied, for example, for any bounded kernel.

The RKHS \mathcal{H}_a can be decomposed as a sum of two orthogonal sub-RKHS,

$$\mathcal{H}_a = \mathcal{H}_{0a} \oplus^\perp \mathcal{H}_{1a},$$

where \mathcal{H}_{0a} is the RKHS of zero mean functions,

$$\mathcal{H}_{0a} = \left\{ f_a \in \mathcal{H}_a, E_{X_a}(f_a(X_a)) = 0 \right\},$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
66 in the regression model with non-Gaussian and non-bounded error**

and \mathcal{H}_{1a} is the RKHS of constant functions,

$$\mathcal{H}_{1a} = \left\{ f_a \in \mathcal{H}_a, f_a(X_a) = C \right\}.$$

The kernel k_{0a} associated with the RKHS \mathcal{H}_{0a} is defined as follows:

$$k_{0a}(X_a, X'_a) = k_a(X_a, X'_a) - \frac{E_{U \sim P_a}(k_a(X_a, U))E_{U \sim P_a}(k_a(X'_a, U))}{E_{(U, V) \sim P_a \otimes P_a} k_a(U, V)}.$$

Let $k_v(X_v, X'_v) = \prod_{a \in v} k_{0a}(X_a, X'_a)$, then the ANOVA kernel k is defined by:

$$k(X, X') = \prod_{a=1}^d (1 + k_{0a}(X_a, X'_a)) = 1 + \sum_{v \in \mathcal{P}} k_v(X_v, X'_v).$$

For \mathcal{H}_v being the RKHS associated with the kernel k_v , the RKHS associated with the ANOVA kernel k is then defined by:

$$\mathcal{H} = \prod_{a=1}^d \left(\mathbb{1} \oplus \mathcal{H}_{0a} \right) = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v,$$

where \perp denotes the L^2 inner product.

According to this construction, any function $f \in \mathcal{H}$ satisfies the following decomposition,

$$f(X) = \langle f, k(X, \cdot) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v).$$

which is the Hoeffding decomposition of f .

For more background on the RKHS spaces see [Aronszajn \(1950\)](#), [Saitoh \(1988\)](#), [Berlinet and Thomas-Agnan \(2003\)](#).

3.2.2 Approximating the Hoeffding decomposition of m

Let $f^* \in \mathcal{H}$ be defined as follows:

$$f^* = \arg \min_{f \in \mathcal{H}} \|m - f\|_2^2 = \arg \min_{f \in \mathcal{H}} E_X (m(X) - f(X))^2.$$

The function $f^* = f_0^* + \sum_{v \in \mathcal{P}} f_v^*$, is the approximation of m on the RKHS \mathcal{H} , and its Hoeffding decomposition is an approximation of the Hoeffding decomposition of m . Therefore, according to Equation (3.3), for all $v \in \mathcal{P}$, each function f_v^* approximates the function m_v .

The number of functions f_v^* is related to the cardinality of \mathcal{P} , i.e. $2^d - 1$, that may be huge. The idea is to calculate a sparse estimator of f^* as an estimator of m . To do so, the ridge group sparse procedure as proposed by [Huet and Taupin \(2017\)](#) is used that we recall in the following.

3.2.3 Ridge group sparse procedure and associated estimator

Let n be the number of observations. For all $v \in \mathcal{P}$, let X_v be the matrix of variables corresponding to the v -th group, i.e.

$$X_v = (X_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}.$$

For any $f \in \mathcal{H}$ such that $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, and for some tuning parameters $\gamma_v, \mu_v, v \in \mathcal{P}$, the ridge group sparse criterion is defined as follows:

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \left(Y_i - f_0 - \sum_{v \in \mathcal{P}} f_v(X_{vi}) \right)^2 + \sum_{v \in \mathcal{P}} \gamma_v \|f_v\|_n + \sum_{v \in \mathcal{P}} \mu_v \|f_v\|_{\mathcal{H}_v},$$

where $\|f_v\|_n$ is the empirical L^2 -norm of f_v defined by the sample $\{X_{vi}\}_{i=1}^n$ as

$$\|f_v\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_v^2(X_{vi}).$$

The penalty function in the criterion $\mathcal{L}(f)$ is the sum of the Hilbert norm and the empirical norm, which allows to select few terms in the additive decomposition of f over sets $v \in \mathcal{P}$. Moreover, the Hilbert norm favours the smoothness of the estimated $f_v, v \in \mathcal{P}$.

Let us define the set of functions,

$$\mathcal{F} = \left\{ f : f = f_0 + \sum_{v \in \mathcal{P}} f_v, \text{ with } f_v \in \mathcal{H}_v, \text{ and } \|f_v\|_{\mathcal{H}_v} \leq r_v, r_v > 0 \right\}. \quad (3.6)$$

Then the RKHS ridge group sparse estimator is defined by,

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f). \quad (3.7)$$

3.3 Risk upper bounds

In this Section, the upper bounds of the empirical L^2 risk and the L^2 risk of the RKHS ridge group sparse estimator are presented in Theorem 3.3.1 and Corollary 3.3.1, respectively. Before stating these results, let us introduce some notation and assumptions that are needed in the rest of this Chapter.

For a function $f \in \mathcal{H}$, let S_f be its support,

$$S_f = \{v \in \mathcal{P} : f_v \neq 0\}. \quad (3.8)$$

The RKHS construction as described in Section 3.2.1 insures that the following properties are satisfied:

- ✓ for all $v \in \mathcal{P}$, the functions $f_v \in \mathcal{H}_v$ are centered and are square-integrable, i.e.

$$E_X(f_v(X_v)) = 0 \text{ and } E_X(f_v^2(X_v)) < \infty,$$

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
68 in the regression model with non-Gaussian and non-bounded error

✓ for all $v, v' \in \mathcal{P}$, $v \neq v'$, the functions $f_v \in \mathcal{H}_v$ and $f_{v'} \in \mathcal{H}_{v'}$ are orthogonal with respect to $L^2(\mathcal{X}, P_X)$, i.e.

$$E_X(f_v(X_v)f_{v'}(X_{v'})) = 0.$$

We assume moreover that,

✓ for all $v \in \mathcal{P}$, the functions $f_v \in \mathcal{H}_v$ are uniformly bounded, i.e.

$$\exists R > 0 \text{ such that } \|f_v\|_\infty = \sup_{X_v} |f_v(X_v)| \leq R.$$

Each kernel k_v , $v \in \mathcal{P}$ is associated with an integral operator T_{k_v} from $L^2(\mathcal{X}_v, P_v)$ to $L^2(\mathcal{X}_v, P_v)$ defined by:

$$\forall f \in L^2(\mathcal{X}_v, P_v), T_{k_v}(f) = \int_{\mathcal{X}_v} k_v(\cdot, t)f(t)dP_v(t).$$

For each $v \in \mathcal{P}$, let $\omega_{v,1} \geq \omega_{v,2} \geq \dots \geq 0$ be the eigenvalues of the integral operator T_{k_v} (see Equation (3.19)). Let us define the function $Q_{n,v}(t)$ for some positive t as follows:

$$Q_{n,v}(t) = \sqrt{\frac{5}{n} \sum_{\ell \geq 1} \min(t^2, \omega_{v,\ell})}, \quad (3.9)$$

and for some $\Delta > 0$ let $\nu_{n,v}$ be defined by:

$$\nu_{n,v} = \inf_t \left\{ Q_{n,v}(t) \leq \Delta t^2 \right\}. \quad (3.10)$$

For each $v \in \mathcal{P}$, $\nu_{n,v}$ refers to the minimax optimal rate for $L^2(\mathcal{X}, P_X)$ -estimation in the RKHS \mathcal{H}_v (Mendelson (2002)).

Remark 3.3.1 *The rate $\nu_{n,v}$, $v \in \mathcal{P}$, depends on the regularity of the RKHS via the decreasing rate of the eigenvalues $\{\omega_{v,\ell}\}_{\ell=1}^\infty$. When RKHS is of high regularity, i.e. when the eigenvalues $\{\omega_{v,\ell}\}_{\ell=1}^\infty$ decrease quickly, then the rate $\nu_{n,v}$, $v \in \mathcal{P}$ will be close to the parametric rate of convergence (see Section 3.3.1).*

The choice of tuning parameters in the criterion $\mathcal{L}(f)$ is specified in terms of the following quantity:

$$\lambda_{n,v} = \max \left(\nu_{n,v}, \sqrt{\frac{d}{n}} \right). \quad (3.11)$$

Theorem 3.3.1 *Consider the regression model defined at Equation (3.1) with $\sigma = 1$. Let $\{(Y_i, X_i)\}_{i=1}^n$ be a n -sample with the same law as (Y, X) , and let $\{\varepsilon_i\}_{i=1}^n$ be the random errors that are independent and identically distributed (i.i.d.) like ε . Let also \hat{f} be defined by (3.7) with $r_v = 1$ in (3.6), and let the tuning parameters μ_v 's and γ_v 's be chosen as follows:*

For some constant $C_1 > 10 + 4\Delta$,

$$\forall v \in \mathcal{P}, \mu_v = C_1 \lambda_{n,v}^2, \gamma_v = C_1 \lambda_{n,v}. \quad (3.12)$$

If there exists positive constants C_2, C_3 , and $0 < \beta < 1/\alpha$ such that the following assumptions are satisfied:

$$\forall v \in \mathcal{P}, n\lambda_{n,v}^2 \geq -C_2 \log \lambda_{n,v}, \quad (3.13)$$

and

$$\forall f \in \mathcal{F}, \sum_{v \in S_f} \lambda_{n,v}^2 \leq C_3 n^{2\beta-1}, \quad (3.14)$$

then, there exists $0 < \eta < 1$ depending on constants $\{C_i\}_{i=1}^3$, β , and n (η tends to 0 as n increases), such that with probability greater than $1 - \eta$, we have for some constant C ,

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\}. \quad (3.15)$$

Let us now comment on the theorem.

Remark 3.3.2 Let f' be the function in \mathcal{F} such that the infimum of the right hand side of the inequality (3.15) is realized. The term $\|m - f'\|_n^2$ is the usual bias term. It quantifies both the approximation properties of the RKHS \mathcal{H} , and the bias-variance trade-off.

Remark 3.3.3 This result is similar to the one obtained in the Gaussian regression model at the cost of the additional Assumption (3.14). This assumption allows to obtain the same rate of convergence for the RKHS ridge group sparse estimator as in the Gaussian regression model (see [Huet and Taupin \(2017\)](#)). However, it implies some restrictions on the regularity of the RKHS \mathcal{H} . Indeed, as for all $v \in \mathcal{P}$, $\lambda_{n,v} \geq \nu_{n,v}$ (see Equation (3.11)), it follows that $\sum_{v \in S_f} \nu_{n,v}^2 \leq C_3 n^{2\beta-1}$, which implies some restrictions on the regularity of the RKHS: if β is small, which will be the case if α is large, then the RKHS should be of high regularity.

Remark 3.3.4 By Equation (3.11), we also have that for all $v \in \mathcal{P}$, $\lambda_{n,v} \geq \sqrt{d/n}$. This assumption allows to control the probability of the $|\mathcal{P}|$ events (see Equation (3.48)), where $\log(|\mathcal{P}|)$ is of order d .

Remark 3.3.5 The result in Theorem 3.3.1 can be generalized to the case where $\sigma \neq 1$ in Equation (3.1), and where $r_v \neq 1$ in (3.6).

Let \hat{g} be defined as follows:

$$\hat{g} = \arg \min_{g \in \mathcal{F}'} \left\{ \left\| \frac{Y}{\sigma} - g \right\|_n^2 + \frac{1}{\sigma} \sum_v \gamma_v \|g_v\|_n + \frac{1}{\sigma} \sum_v \mu_v \|g_v\|_{\mathcal{H}_v} \right\}, \quad (3.16)$$

with

$$\mathcal{F}' = \left\{ g : g = g_0 + \sum_v g_v, \text{ with } g_v \in \mathcal{H}_v, \text{ and } \|g_v\|_{\mathcal{H}_v} \leq \frac{r_v}{\sigma} \right\}. \quad (3.17)$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
70 in the regression model with non-Gaussian and non-bounded error**

We have $\hat{f} = \sigma \hat{g}$ for \hat{f} being defined by (3.7).

For all $u > 0$, let \mathcal{H}_v^u be the RKHS associated with the kernel uk_v . If $u = r_v^2/\sigma^2$, then

$$\hat{g} = \arg \min_{g \in \mathcal{F}''} \left\{ \left\| \frac{Y}{\sigma} - g \right\|_n^2 + \frac{1}{\sigma} \sum_v \gamma_v \|g_v\|_n + \frac{1}{\sigma^2} \sum_v \mu_v r_v \|g_v\|_{\mathcal{H}_v^u} \right\}.$$

where

$$\mathcal{F}'' = \left\{ g : g = g_0 + \sum_v g_v, \text{ with } g_v \in \mathcal{H}_v^u, \text{ and } \|g_v\|_{\mathcal{H}_v^u} \leq 1 \right\}.$$

We apply Theorem 3.3.1 with Y/σ and m/σ in place of Y and m , to \hat{g} defined as above.

Let

$$Q_{n,v}^u(t) = \sqrt{\frac{5}{n} \sum_{\ell \geq 1} \min(t^2, u\omega_{v,\ell})},$$

and for $\Delta' > 0$, let

$$\nu_{n,v}^u(\Delta') = \inf_t \left\{ Q_{n,v}^u(t) \leq \Delta' t^2 \right\}.$$

Let also

$$\lambda_{n,v}^u = \max \left(\nu_{n,v}^u, \sqrt{\frac{d}{n}} \right).$$

For some constant $C_1 > 10 + \Delta'$, take

$$\frac{\mu_v r_v}{\sigma^2} = C_1 \left(\lambda_{n,v}^u \right)^2, \quad \frac{\gamma_v}{\sigma} = C_1 \lambda_{n,v}^u.$$

Then, for S_g being defined as follows

$$S_g = \{v \in \mathcal{P} : g_v \neq 0\}, \tag{3.18}$$

we have

$$\left\| \frac{m}{\sigma} - \hat{g} \right\|_n^2 \leq C \inf_{g \in \mathcal{F}''} \left\{ \left\| \frac{m}{\sigma} - g \right\|_n^2 + \frac{1}{\sigma^2} \sum_{v \in S_g} (\mu_v r_v + \gamma_v^2) \right\},$$

or, multiplying both sides by σ^2 , and taking $u = r_v^2/\sigma^2$,

$$\|m - \sigma \hat{g}\|_n^2 \leq C \inf_{g \in \mathcal{F}'} \left\{ \|m - \sigma g\|_n^2 + \sum_{v \in S_g} (\mu_v r_v + \gamma_v^2) \right\}.$$

Corollary 3.3.1 Under the same assumptions as Theorem 3.3.1, we have with high probability for some constant C' that,

$$\|m - \hat{f}\|_2^2 \leq C' \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \|m - f\|_2^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) \right\}.$$

Remark 3.3.6 *The result in Corollary 3.3.1 can be generalized to the case where $\sigma \neq 1$ in Equation (3.1), and where $r_v \neq 1$ in (3.6). It suffices to apply Corollary 3.3.1 with Y/σ and m/σ in place of Y and m , to \hat{g} as defined in Equation (3.16). Then, with similar demonstration as in Remark 3.3.5 we obtain,*

$$\|m - \sigma \hat{g}\|_2^2 \leq C' \inf_{g \in \mathcal{F}'} \left\{ \|m - \sigma g\|_n^2 + \|m - \sigma g\|_2^2 + \sum_{v \in S_g} (\mu_v r_v + \gamma_v^2) \right\},$$

where \mathcal{F}' and S_g are defined in Equations (3.17) and (3.18), respectively.

3.3.1 Rate of convergence

Corollary 3.3.2 *Under the same assumptions as Theorem 3.3.1, we have*

$$\|m - \hat{f}\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \left(\sum_{v \in S_f} \nu_{n,v}^2 + \frac{d|S_f|}{n} \right) \right\}.$$

This Corollary highlights that the upper bound is relevant when the infimum is reached for functions f that have a sparse decomposition in \mathcal{H} , i.e. $|S_f|$ is small, and when d is small face to n . When d is large, the decomposition of functions in \mathcal{H} should be limited to interactions of a limited order, so that the number of elements in the estimated meta-model is of order smaller than d^r for some small r , say $r = 2$ for example. In such a case, the cardinality of \mathcal{P} will be smaller than d^2 . As we mentioned in Remark 3.3.4, the assumption $\lambda_{n,v} \geq \sqrt{d/n}$ is needed to control the value $\log(|\mathcal{P}|)$, which will be now smaller than $2 \log(d)$. Therefore, the value d in the definition of $\lambda_{n,v}$ (see Equation (3.11)) as well as the term $d|S_f|/n$ in the infimum above will be replaced by $2 \log(d)$ and $2 \log(d)|S_f|/n$, respectively.

Let us discuss the rate of convergence given by $\sum_{v \in S_f} \nu_{n,v}^2$. For the sake of simplicity we consider the case where the variables X_1, \dots, X_d have the same distribution P_1 on $\mathcal{X}_1 \subset \mathbb{R}$, and where the unidimensional kernels k_{0a} are all identical, such that $k_v(X_v, X'_v) = \prod_{a \in v} k_0(X_a, X'_a)$. The kernel k_0 admits an eigen expansion given by

$$k_0(X_a, X'_a) = \sum_{\ell_a \geq 1} \omega_{0,\ell_a} \phi_{\ell_a}(X_a) \phi_{\ell_a}(X'_a),$$

where the eigenvalues $\{\omega_{0,\ell_a}\}_{\ell_a=1}^\infty$ are non-negative and ranged in the decreasing order, and where the $\{\phi_{\ell_a}\}_{\ell_a=1}^\infty$ are the associated eigenfunctions, orthonormal with respect to $L^2(\mathcal{X}_1, P_1)$. Therefore, the kernel k_v admits the following expansion,

$$k_v(X_v, X'_v) = \sum_{\ell=(\ell_1 \dots \ell_{|v|})} \underbrace{\prod_{a=1}^{|v|} \omega_{0,\ell_a}}_{\omega_{v,\ell}} \underbrace{\prod_{a=1}^{|v|} \phi_{\ell_a}(X_a)}_{\phi_{v,\ell}(X_v)} \underbrace{\prod_{a=1}^{|v|} \phi_{\ell_a}(X'_a)}_{\phi_{v,\ell}(X'_v)}. \quad (3.19)$$

Consider the case where the eigenvalues $\{\omega_{0,\ell_a}\}_{\ell_a=1}^\infty$ are decreasing at a rate $\ell_a^{-2\alpha'}$ for some $\alpha' > 1/2$, i.e. the $\omega_{0,\ell}$ are of order $\ell^{-2\alpha'} = (\prod_{a=1}^{|v|} \ell_a)^{-2\alpha'}$. It is shown in

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
72 in the regression model with non-Gaussian and non-bounded error**

Section 8.3. of [Huet and Taupin \(2017\)](#), that

$$\nu_{n,v} \propto n^{-\frac{\alpha'}{(2\alpha'+1)}} (\log n)^{\gamma'},$$

where the rate $\nu_{n,v}$ is defined at Equation (3.10) and

$$\gamma' \geq (|v| - 1) \frac{\alpha'}{(2\alpha' - 1)}.$$

For all $f \in \mathcal{F}$ we have then,

$$\sum_{v \in S_f} \nu_{n,v}^2 \propto |S_f| n^{-\frac{2\alpha'}{(2\alpha'+1)}} (\log n)^{2\gamma'}.$$

Note that in this particular case, the rate of convergence depends on $|v|$ through the logarithmic term $(\log n)^{2\gamma'}$, and that up to this logarithmic term the rate of convergence has the same order than the usual non-parametric rate for unidimensional functions. It follows that the RKHS space \mathcal{H} should be chosen such that the unknown function m is well approximated by sparse functions in \mathcal{H} with low order of interactions.

Besides, the rate $\nu_{n,v}$ should satisfy assumption (3.14),

$$\sum_{v \in S_f} \nu_{n,v}^2 \leq C_3 n^{2\beta-1},$$

which holds if

$$\alpha' > \frac{1 - 2\beta}{4\beta} > \frac{\alpha - 2}{4}. \quad (3.20)$$

This shows that for the large values of α the assumption (3.14) implies some restrictions on the regularity of the RKHS chosen: If $\alpha < 4$, then all α' greater than $1/2$ satisfy Equation (3.20), since $(\alpha - 2)/4 < 1/2$. If $\alpha \geq 4$, then we have $\alpha' > (\alpha - 2)/4 > 1/2$. As α increases, i.e. β decreases (recall that $0 < \beta < 1/\alpha$), and assumption (3.14) implies that the RKHS chosen should be of high regularity.

3.4 Main arguments of the proof of Theorem 3.3.1 and motivation for the choice π_α

The proof of Theorem 3.3.1 starts in the same way as the proof of Theorem 2.1. in [Huet and Taupin \(2017\)](#) where they considered the Gaussian regression model. However, it differs in two essential points:

1. Sudakov type minoration,
2. Concentration inequality.

In the following Section, we give a sketch of the proof of Theorem 3.3.1, we highlight the two points above that differs the proof from the proof in the Gaussian regression model, and we provide a detailed comparison to the related works. In Section 3.4.2 we give a brief introduction to the Sudakov type minoration context, we explain the motivation for choosing densities $\pi_\alpha \in \mathcal{D}$ defined in Equation (3.2), and we state in Corollary 3.4.1 the appropriate Sudakov minoration used in the proof of Theorem 3.3.1. In Section 3.4.3 we present the concentration inequality context, and we state in Corollary 3.4.2 the appropriate concentration inequality used in the proof of Theorem 3.3.1.

3.4.1 Sketch of the proof

We give here a sketch of the proof of Theorem 3.3.1, and we postpone to Section 3.5 for complete statements. We begin by introducing some notation.

We denote by C constants that vary from an equation to the other. For $v \in \mathcal{P}$, and for a function $\phi : \mathbb{R}^{|v|} \mapsto \mathbb{R}$, we denote by $V_{n,\varepsilon}$ the empirical process defined as,

$$V_{n,\varepsilon}(\phi) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i \phi(X_{v,i}). \quad (3.21)$$

For all $v \in \mathcal{P}$, let \mathcal{H}_v be the RKHS associated with the reproducing kernel k_v . For any function $g_v \in \mathcal{H}_v$, $v \in \mathcal{P}$, and $V_{n,\varepsilon}$ being defined in Equation (3.21), we consider two following processes,

$$W_{n,2,v}(t) = \sup \left\{ |V_{n,\varepsilon}(g_v)|, \|g_v\|_{\mathcal{H}_v} \leq 2, \|g_v\|_2 \leq t \right\}, \quad (3.22)$$

$$W_{n,n,v}(t) = \sup \left\{ |V_{n,\varepsilon}(g_v)|, \|g_v\|_{\mathcal{H}_v} \leq 2, \|g_v\|_n \leq t \right\}. \quad (3.23)$$

Starting from the definition of \hat{f} , some simple calculations give that for all $f \in \mathcal{F}$,

$$\begin{aligned} C \|m - \hat{f}\|_n^2 &\leq \|m - f\|_n^2 + |V_{n,\varepsilon}(\hat{f} - f)| + \sum_{v \in S_f} [\gamma_v \|\hat{f}_v - f_v\|_n + \mu_v \|\hat{f}_v - f_v\|_{\mathcal{H}_v}] \\ &\quad - \sum_{v \notin S_f} [\mu_v \|\hat{f}_v\|_{\mathcal{H}_v} + \gamma_v \|\hat{f}_v\|_n], \\ &\leq \|m - f\|_n^2 + |V_{n,\varepsilon}(\hat{f} - f)| + \sum_{v \in S_f} [\gamma_v \|\hat{f}_v - f_v\|_n + \mu_v \|\hat{f}_v - f_v\|_{\mathcal{H}_v}]. \end{aligned}$$

If we set $g = \hat{f} - f$, then $g \in \mathcal{H}$, $g = g_0 + \sum_v g_v$, with $g_v = \hat{f}_v - f_v$, and for each v , $\|g_v\|_{\mathcal{H}_v} \leq 2$.

The main problem is now to control the empirical process $V_{n,\varepsilon}$. For each v , letting $\lambda_{n,v}$ as in (3.11), we state (see Lemma 3.5.1, page 85) that, with high probability,

$$|V_{n,\varepsilon}(g_v)| \leq C \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + C \lambda_{n,v} \|g_v\|_n. \quad (3.24)$$

One of the key points in the proof of Lemma 3.5.1 is to find an upper bound for the two following quantities:

$$|W_{n,n,v}(t) - E_\varepsilon(W_{n,n,v}(t))|, \text{ and } |W_{n,2,v}(t) - E_\varepsilon(W_{n,2,v}(t))|. \quad (3.25)$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
74 in the regression model with non-Gaussian and non-bounded error**

In the Gaussian regression model, one use the isoperimetric inequality for Gaussian processes in [Massart and Picard \(2007\)](#).

When dealing with errors that are not distributed as a Gaussian distribution, different tools are needed to obtain the upper bounds for the quantities in Equation (3.25) (see Section 3.4.3 for a complete discussion of this point of the proof). Let us continue the sketch of the proof before coming back to this point.

If for all v , μ_v and γ_v satisfying Equation (3.12), by using Equation (3.24) we deduce that with high probability,

$$C\|m - \hat{f}\|_n^2 \leq \|m - f\|_n^2 + \sum_{v \in S_f} [\gamma_v \|g_v\|_n + \mu_v \|g_v\|_{\mathcal{H}_v}] + \sum_{v \notin S_f} [\gamma_v \|\hat{f}_v\|_n + \mu_v \|\hat{f}_v\|_{\mathcal{H}_v}].$$

Besides, we can express the decomposability property of the penalty as follows (see lemma 3.5.2, page 85):

over the set where the empirical process is controlled as stated above, we have with high probability,

$$\sum_{v \notin S_f} [\gamma_v \|\hat{f}_v\|_n + \mu_v \|\hat{f}_v\|_{\mathcal{H}_v}] \leq C \sum_{v \in S_f} [\gamma_v \|g_v\|_n + \mu_v \|g_v\|_{\mathcal{H}_v}].$$

Putting the things together, and using that $\|g_v\|_{\mathcal{H}_v} \leq 2$, we obtain the following upper bound:

$$C\|m - \hat{f}\|_n^2 \leq \|m - f\|_n^2 + \sum_{v \in S_f} [\mu_v + \gamma_v \|g_v\|_n].$$

The last important step consists in comparing $\sum_{v \in S_f} \|g_v\|_n$ to $\|\sum_{v \in S_f} g_v\|_n$. To do so, we show first (see lemma 3.5.3 page 86) that for all $v \in \mathcal{P}$, with high probability,

$$\|g_v\|_n \leq 2\|g_v\|_2 + \gamma_v.$$

Using inequality above and that for all positive K , $2ab \leq (1/K)a^2 + Kb^2$ we obtain,

$$\begin{aligned} C\|m - \hat{f}\|_n^2 &\leq \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) + \sum_{v \in S_f} \|g_v\|_2^2 \\ &\leq \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) + \sum_{v \in \mathcal{P}} \|g_v\|_2^2. \end{aligned}$$

Then we use the orthogonality assumption between the spaces \mathcal{H}_v ,

$$\sum_{v \in \mathcal{P}} \|g_v\|_2^2 = \left\| \sum_{v \in \mathcal{P}} g_v \right\|_2^2 = \|g\|_2^2,$$

which allows us to obtain the following result:

$$C\|m - \hat{f}\|_n^2 \leq \|m - f\|_n^2 + \sum_{v \in S_f} (\mu_v + \gamma_v^2) + \|\hat{f} - f\|_2^2.$$

It remains now to consider different cases according to the rankings of $\|\hat{f} - f\|_2^2$ and $\|\hat{f} - f\|_n^2$ to get the result of Theorem 3.3.1.

If $\|\widehat{f} - f\|_2 \leq \|\widehat{f} - f\|_n$ the result is obtained by a simple rearrangement of the terms.

If $\|\widehat{f} - f\|_2 \geq \|\widehat{f} - f\|_n$, under some suitable assumptions it is shown (see Lemma 3.5.4 page 87) that with high probability we have

$$\|\widehat{f} - f\|_2 \leq \sqrt{2}\|\widehat{f} - f\|_n.$$

One of the steps to prove the inequality above is to lower bound the expectation of the supremum of the empirical process, i.e. $\mathbb{E}_\varepsilon \sup_g |V_{n,\varepsilon}(g)|$ by a function of the covering number of the functional class under study, say \mathcal{G} . In order to solve this step in the Gaussian regression model one may use the Sudakov minoration in [Pisier \(1989\)](#), for which the minoration is obtained thanks to the Slepian's Lemma. The Slepian's Lemma is specific to the Gaussian setting, and it does not hold when dealing with errors that are not distributed as a centered Gaussian distribution.

In the regression model (see Equation (3.1)) with error ε that is distributed with density proportional to $\pi_\alpha \in \mathcal{D}$, the proof of the upper bound stated in Theorem 3.3.1, needs two following mathematical tools:

- Point 1. a Sudakov type minoration to link the covering number on a class \mathcal{G} to the expectation of the supremum of the empirical process over this class \mathcal{G} , $\mathbb{E}_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|$, and conclude Lemma 3.5.4,
- Point 2. a concentration inequality to bound the quantities defined in Equation (3.25) which leads to bound the empirical process $V_{n,\varepsilon}$ and conclude Lemma 3.5.1.

The Point 1. is solved using a Sudakov type minoration which is a consequence of the result obtained by [Talagrand \(1994\)](#). More precisely, it can be shown (see Corollary 3.4.1 page 80) that for $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ being i.i.d. random variables distributed with density $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)), and for all $\delta > 0$, we have,

$$\begin{aligned} \frac{1}{K} \log N(\delta, \mathcal{G}, \|\cdot\|) &\leq \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^2 \mathbf{1}_{[2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|, \infty)}(\delta) \\ &\quad + \left(\frac{2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|}{\delta} \right)^\alpha \mathbf{1}_{(0, 2nE_\varepsilon \sup_{g \in \mathcal{G}} |V_{n,\varepsilon}(g)|]}(\delta), \end{aligned} \quad (3.26)$$

where K is a constant that depends on α only, $\|\cdot\|$ is the Euclidean norm, $N(\delta, \mathcal{G}, \|\cdot\|)$ is the δ -covering number of the metric space $(\mathcal{G}, \|\cdot\|)$, and $1_A : \mathcal{A} \rightarrow \{0, 1\}$ is the indicator function of $A \subset \mathcal{A}$, i.e.

$$1_A(a) = \begin{cases} 1 & \text{if } a \in A, \\ 0 & \text{if } a \notin A. \end{cases}$$

The proof of Lemma 3.5.4 proceeds using Equation (3.24) and is concluded under the Hypothesis (3.12) and (3.14).

The Point 2. is solved using a concentration inequality (see Corollary 3.4.2 page 83) which is a consequence of the result obtained by [Shu and Strzelecki \(2017\)](#). □

The appropriate results to solve Point 1. and Point 2. are stated in Corollary 3.4.1 in Section 3.4.2.2 and Corollary 3.4.2 in Section 3.4.3.2, respectively.

3.4.1.1 Comparison with related works

- ✓ Meier et al. (2009) considered a least-squares criterion penalized by a penalty function similar to the one we consider in our work. Their estimator of the unknown function m has an univariate additive decomposition, i.e. decomposition (3.3) limited to the main effects.

They used a *compatibility condition* to compare the sum of the empirical L^2 -norm of the univariate functions to the empirical L^2 -norm of the sum of the univariate functions. More precisely,

Let $S^* = \{a \in \{1, \dots, d\}, \|f_a\|_n \neq 0\}$, then for $C(f_a)$ being a term depending on the functions f_a , $a \in S^*$,

$$\sum_{a \in S^*} \|f_a\|_n^2 \leq \left\| \sum_{a=1}^d f_a \right\|_n^2 + C(f_a).$$

The control of the Empirical process is done in their Lemma 1. This lemma is proved using Lemma 8.4 in van de Geer et al. (2000), for which the errors should have sub-Gaussian tails, i.e.

$$\max_i E \left(\exp \left(\frac{\varepsilon_i^2}{C_1} \right) \right) \leq C_2,$$

where C_1 and C_2 are constants.

Afterwards, it was shown by Raskutti et al. (2012) (see Section 3.4. of their paper) that the convergence rate of this estimator is sub-optimal.

- ✓ Koltchinskii and Yuan (2010) considered a large class of loss functions, called *loss functions of quadratic type*, which satisfies the boundedness conditions. More precisely, for l being a loss function, they assume that $l(Y, \cdot)$ is uniformly bounded from above by a numerical constant. So for a given distribution of the observations Y , there may exists a loss function that belongs to the class of the *loss functions of quadratic type* (see Section 2.1. of their paper for some examples).

They consider the input variables X that may be not independent, and they do not assume that there is orthogonality between their RKHS, therefore $\|\sum_v f_v\|_2 \neq \sum_v \|f_v\|_2$. Instead, in their Section 2.2., they introduce some geometric characteristics related to the degree of *dependence* of their RKHS, which insures *almost* orthogonality between these spaces.

The control of the empirical process is done in their Lemma 9. This lemma is proved under the global boundedness condition and the assumptions of the *loss functions of quadratic type*.

We consider the quadratic loss function to obtain an estimator of the function m in the regression model defined in Equation (3.1), with error ε that is non-bounded. This case is not included in the class of the *loss functions of*

quadratic type. We do not impose the global boundedness condition. Instead, we assume that for all $v \in \mathcal{P}$ the functions f_v are uniformly bounded. More precisely, the quantity $\sup_{X \in \mathcal{X}} |f_v(X)|$ is bounded from above by a constant. This assumption is easily satisfied as soon as the kernel k_v is bounded on the compact set \mathcal{X} ,

$$\sup_{X \in \mathcal{X}} |f_v(X)| \leq \sup_{X \in \mathcal{X}} \sqrt{k_v(X_v, X_v)} \|f_v\|_{\mathcal{H}_v}.$$

For a detailed discussion on this subject, we refer to the paper by [Raskutti et al. \(2012\)](#).

- ✓ In the Gaussian regression model,
 - [Raskutti et al. \(2012\)](#) assumed that the unknown function m has a sparse univariate decomposition, where each component in its decomposition lies in a RKHS. They obtained an estimator for m , based on a ridge group sparse type procedure. They established upper and lower bounds on the risk in the L^2 -norm and upper bound on the risk in the empirical L^2 -norm.
 - [Huet and Taupin \(2017\)](#) assumed that the unknown function m admits a Hoeffding decomposition involving the main effects and interactions. They obtained a RKHS ridge group sparse estimator of a meta-model that approximates the Hoeffding decomposition of m . They established upper bounds on the risk in the L^2 -norm and the empirical L^2 -norm.

[Raskutti et al. \(2012\)](#) and [Huet and Taupin \(2017\)](#) do not assume global boundedness condition. Instead, they assume that for all $v \in \mathcal{P}$ the functions f_v are uniformly bounded. The proof of their results relies on the empirical Gaussian process methods such as Sudakov minoration [Pisier \(1989\)](#) and concentration inequalities for Gaussian processes.

As we are not in the Gaussian regression model, these methods could not be used in our work. We require new tools that we describe in details in the two next Sections.

3.4.2 Sudakov minoration

In the following Section, we recall the definition of the covering numbers, the statement of the classical Sudakov minoration, which is specific to the Gaussian process, and the generalized Sudakov minoration known also as the Sudakov minoration principal, which could be applied to some other processes. In Section 3.4.2.2 we state the appropriate Sudakov type minoration to the process associated with the random variables that are distributed with density $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)) in Corollary 3.4.1.

3.4.2.1 Introduction

Let T be a set of square-integrable functions, i.e. $T \subset L^2$, and $\|\cdot\|$ be the Euclidean norm. For any $\delta > 0$, we denote by $C(\delta, T, \|\cdot\|)$ the δ -covering set of the metric space $(T, \|\cdot\|)$:

$$C(\delta, T, \|\cdot\|) = \left\{ f^1, \dots, f^N : \forall f \in T, \exists k \in \{1, \dots, N\} \text{ such that } \|f - f^k\| \leq \delta \right\}.$$

The δ -covering number of $(T, \|\cdot\|)$, denoted $N(\delta, T, \|\cdot\|)$, is the cardinal of the smallest covering set. A proper covering restricts the covering to use only elements in the set T . It can be shown that the covering numbers and the proper covering numbers are related by the following inequality:

$$N(\delta, T, \|\cdot\|) \leq N_{\text{proper}}(\delta, T, \|\cdot\|) \leq N\left(\frac{\delta}{2}, T, \|\cdot\|\right). \quad (3.27)$$

Consider a random variable Z such that $E(Z^2) < \infty$, and consider an i.i.d. sequence $\{Z_i\}_{i=1}^n$ distributed like Z . To each $t = (t_1, \dots, t_n)$ of $T \subset L^2$ one can associate the process $V_t = \sum_{i=1}^n Z_i t_i$, $t \in T$.

In order to link the covering number on a class T , i.e. $N(\delta, T, \|\cdot\|)$, to the expectation of the supremum of the process $V_t = \sum_{i=1}^n Z_i t_i$ in the Gaussian setting, the classical Sudakov minoration could be used (Pisier (1989)):

$$\frac{1}{K} \log N(\delta, T, \|\cdot\|) \leq \left(\frac{n E_Z \sup_{t \in T} \sum_{i=1}^n Z_i t_i}{\delta} \right)^2. \quad (3.28)$$

When dealing with the processes $V_t = \sum_{i=1}^n Z_i t_i$, $t \in T$ associated with the random variables $\{Z_i\}_{i=1}^n$ that are not Gaussian, a generalized Sudakov minoration, known also as the Sudakov minoration principal, could be used to lower bound the value $E_Z \sup_{t \in T} \sum_{i=1}^n Z_i t_i$. Let us recall this inequality.

Definition 3.4.1 (Definition 1.1. in Latała (2014)) *Let $Z = (Z_1, \dots, Z_n)$ be a random vector in \mathbb{R}^n . We say that Z satisfies the L_p -Sudakov minoration principle with a constant $K' > 0$, $SMP_p(K')$, if for any set $T \subset \mathbb{R}^n$ with $|T| > \exp(p)$ such that*

$$\left(E_Z \sum_{t, s \in T} \left| \sum_{i=1}^n (t_i - s_i) Z_i \right|^p \right)^{1/p} := \left\| \sum_{i=1}^n (t_i - s_i) Z_i \right\|_p \geq \delta, \quad \forall s, t \in T, s \neq t, \quad (3.29)$$

we have

$$K' \delta \leq E_Z \sup_{t, s \in T} \sum_{i=1}^n (s_i - t_i) Z_i.$$

A random vector Z satisfies the Sudakov minoration principle with a constant K' , $SMP(K')$, if it satisfies $SMP_p(K')$ for any $p \geq 1$.

If $\{Z_i\}_{i=1}^n$ are independent symmetric ± 1 random variables or equivalently if the vector $Z = (Z_1, \dots, Z_n)$ is uniformly distributed on the cube $[-1, 1]$ the Sudakov minoration principal with universal K' was proven by Talagrand (1993).

Latała (2014) proved the Sudakov minoration principal for the independent log-concave random variables. A measure on \mathbb{R}^n with the full dimensional support is log-concave if and only if it has a density of the form $\exp(-\phi(x))$, where $\phi : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is convex (Borell (1974)). In the dependent setting the Sudakov minoration principal for the log-concave random variables was proven by Bednorz (2014).

As we are in the independent setting and the densities $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)) are log-concave, the Sudakov minoration obtained by Latała (2014) holds in our context. However, we could not deduce from the result obtained by Latała (2014) the adapted Sudakov type minoration that leads to obtain the *optimal* rate of convergence for our estimator. By *optimal* we mean the same rate of convergence as in the Gaussian regression setting (see Huet and Taupin (2017)). This is the reason why we restricted ourselves to the densities $\pi_\alpha \in \mathcal{D}$ for which there exists a result given by Talagrand (1994).

In the next Section we provide in Corollary 3.4.1 the appropriate Sudakov type minoration for the random variables that are distributed with density $\pi_\alpha \in \mathcal{D}$. This Corollary is a consequence of the result obtained by Talagrand (1994).

3.4.2.2 Sudakov minoration for density π_α

In this Section we state in Corollary 3.4.1 the Sudakov minoration appropriate for the random variables that are distributed with density $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)). This Corollary is a consequence of the Sudakov minoration stated in Theorem 3.1. in Talagrand (1994). We start by introducing some notation that we need in the rest of this Section.

Let us denote by $\tilde{\alpha}$ the conjugate exponent of α , i.e. $1/\alpha + 1/\tilde{\alpha} = 1$. So, for all $\alpha > 2$ we have $1 < \tilde{\alpha} < 2$.

We consider the sets $B_{\tilde{\alpha}}$ and $U_{\tilde{\alpha}}(u)$, $u \geq 0$ defined as follows:

$$B_{\tilde{\alpha}} = \left\{ x \in \mathbb{R}^n : \sum_{k=1}^n |x_k|^{\tilde{\alpha}} \leq 1 \right\}, \quad (3.30)$$

and

$$U_{\tilde{\alpha}}(u) = \left\{ x \in \mathbb{R}^n : \sum_{i=1}^n \eta_{\tilde{\alpha}}(x_i) \leq u, u \geq 0 \right\}, \quad (3.31)$$

where

$$\eta_{\tilde{\alpha}}(x_i) = x_i^2 \mathbf{1}_{[-1,1]}(x_i) + |x_i|^{\tilde{\alpha}} \mathbf{1}_{(-\infty,-1] \cup [1,\infty)}(x_i).$$

For $T \subset L^2$ and $u \geq 0$, let $D(T, U_{\tilde{\alpha}}(u))$ be a covering set of translates of T by $U_{\tilde{\alpha}}(u)$:

$$\begin{aligned} D(T, U_{\tilde{\alpha}}(u)) &= \left\{ f^1, \dots, f^N : \forall f \in T, \exists k \in \{1, \dots, N\} \text{ such that } f - f^k \in U_{\tilde{\alpha}}(u) \right\}, \\ &= \left\{ f^1, \dots, f^N : \forall f \in T, \exists k \in \{1, \dots, N\} \text{ such that } \sum_{i=1}^N \eta_{\tilde{\alpha}}(f_i - f_i^k) \leq u \right\}. \end{aligned}$$

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
80 in the regression model with non-Gaussian and non-bounded error

We denote by $N(T, U_{\tilde{\alpha}}(u))$ the minimum number of translates of $U_{\tilde{\alpha}}(u)$ by elements of T needed to cover T .

Lemma 3.4.1 *For all $\tilde{\alpha} \leq 2$ and $u \geq 0$, it is shown that (Talagrand (1994)):*

$$U_{\tilde{\alpha}}(u) \subset (u^{1/2}B_2 + u^{1/\tilde{\alpha}}B_{\tilde{\alpha}}). \quad (3.32)$$

Remark 3.4.1 *If $\tilde{\alpha} \leq 2$ and $u \geq 0$, then*

$$U_{\tilde{\alpha}}(u) \subset 2 \times \max(u^{1/2}, u^{1/\tilde{\alpha}})B_2.$$

The proof of Remark 3.4.1 is given in Section 3.A.1 page 108.

Theorem 3.4.1 *(Theorem 3.1. in Talagrand (1994)) Let $Z = (Z_1, \dots, Z_n)$ be i.i.d. random variables distributed with density $\pi_{\alpha} \in \mathcal{D}$ defined in Equation (3.2), $U_{\tilde{\alpha}}(u)$, $u \geq 0$ be defined by (3.31) and $T \subset L^2$. Set*

$$M = E_Z \sup_{t \in T} \sum_{i=1}^n t_i Z_i, \quad (3.33)$$

then it is shown that:

$$N(T, U_{\tilde{\alpha}}(M)) \leq \exp(KM), \quad (3.34)$$

where K is a constant that depends on α only.

Remark 3.4.2 *According to Theorem 3.4.1 and Remark 3.4.1 for all $u \geq 0$ we have,*

$$N(2 \times \max(u^{1/2}, u^{1/\tilde{\alpha}}), T, \|\cdot\|) \leq N(T, U_{\tilde{\alpha}}(u)) \leq \exp(Ku). \quad (3.35)$$

To be more precise, since $1 < \tilde{\alpha} < 2$ we have

$$(i) \text{ For } u \leq 1, u^{1/\tilde{\alpha}} \leq u^{1/2} \text{ and } N(2u^{1/2}, T, \|\cdot\|) \leq \exp(Ku).$$

$$(ii) \text{ For } u \geq 1, u^{1/\tilde{\alpha}} \geq u^{1/2} \text{ and } N(2u^{1/\tilde{\alpha}}, T, \|\cdot\|) \leq \exp(Ku).$$

Corollary 3.4.1 *Under the same assumptions as for Theorem 3.4.1 we have for all $\delta > 0$,*

$$\frac{1}{K} \log N(\delta, T, \|\cdot\|) \leq \left(\frac{2M}{\delta}\right)^{\alpha} \mathbf{1}_{(0, 2M]}(\delta) + \left(\frac{2M}{\delta}\right)^2 \mathbf{1}_{[2M, \infty)}(\delta),$$

which is exactly Equation (3.26) with M defined in Equation (3.33).

The proof of Corollary 3.4.1 is given in Section 3.A.2 page 108.

3.4.3 Concentration inequality

We start this Section with a small introduction on the concentration inequalities context in Section 3.4.3.1, and we detail the concentration inequality used in our work in Section 3.4.3.2.

3.4.3.1 Introduction

Let $Z = (Z_1, \dots, Z_n)$ be a random vector in \mathbb{R}^n , and the function ϕ from \mathbb{R}^n to \mathbb{R} be convex and 1-Lipschitz with respect to the Euclidean norm on \mathbb{R}^n , i.e.

$$\|\phi(Z) - \phi(Z')\| \leq \|Z - Z'\|, \quad Z, Z' \in \mathbb{R}^n.$$

We are interested in the concentration inequalities of order two that provide bounds on how $\phi(Z)$ deviates from its expected value. More precisely, for P being the probability measure on \mathbb{R}^n , and for all $u \geq 0$,

$$P\left(|\phi(Z) - E(\phi(Z))| \geq u\right) \leq C_1 \exp\left(-\frac{u^2}{C_2}\right), \quad (3.36)$$

where C_1 , and C_2 are constants.

It was shown by [Ledoux and Talagrand \(1991\)](#) that, if $Z = (Z_1, \dots, Z_n)$ is a centered Gaussian random vector in \mathbb{R}^n , then:

$$P\left(|\phi(Z) - E(\phi(Z))| \geq u\right) \leq 4 \exp\left(-\frac{u^2}{2}\right).$$

This result could be proved using an inequality established by geometric arguments and an induction on the number of coordinates.

After that, an alternative approach to some of Talagrand's inequalities was proposed by [Ledoux \(1997\)](#) based on the log-Sobolev inequalities. He showed that if the probability measure P on $[0, 1]^n$ satisfies the log-Sobolev inequality then it satisfies the concentration inequalities of the form (3.36), i.e. the log-Sobolev inequality implies the deviation inequality.

We say that the probability measure P satisfies the log-Sobolev inequality for a class of functions Ψ with loss function $R : \mathbb{R}^n \rightarrow [0, +\infty)$, if for every $\psi \in \Psi$ we have,

$$\text{Ent}(\exp(\psi)) \leq CE(R(\nabla\psi) \exp(\psi)),$$

where $\nabla\psi$ is the usual gradient of ψ , and $\text{Ent}(\exp(\psi))$ is the usual entropy of $\exp(\psi)$, i.e.

$$\text{Ent}(\exp(\psi)) = E(\psi \exp(\psi)) - E(\exp(\psi)) \log(E(\exp(\psi))).$$

This inequality was first introduced by [Gross \(1975\)](#) with $R(x) = \|x\|^2$, $x \in \mathbb{R}^n$ and Ψ being the class of \mathcal{C}^1 functions. A lot of work has been done with different loss and class of functions, see for example [Bobkov and Ledoux \(1997\)](#), [Gentil et al. \(2005, 2007\)](#).

In the rest of this Chapter, we assume that Ψ is the class of convex functions, and we consider only the quadratic loss $R(x) = \|x\|^2$, $x \in \mathbb{R}^n$. Therefore, the probability measure P satisfies the convex log-Sobolev inequality if,

$$E(\psi \exp(\psi)) - E(\exp(\psi)) \log(E(\exp(\psi))) \leq CE(\|\nabla\psi\|^2 \exp(\psi)). \quad (3.37)$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
82 in the regression model with non-Gaussian and non-bounded error**

Adameczak (2005) found a sufficient condition for a class of probability distributions, denoted $\mathcal{M}(m, \rho^2)$ with $m > 0$ and $\rho \geq 0$, on the real line, to satisfy the convex log-Sobolev inequality. He deduced then the following concentration inequality which is satisfied for all probability distributions belonging to $\mathcal{M}(m, \rho^2)$:

$$P\left(\phi(Z) - E(\phi(Z)) \geq u\right) \leq \exp\left(-\frac{u^2}{4C(m, \rho^2)}\right). \quad (3.38)$$

We show in Lemma 3.4.2 that the probability distributions associated with the densities $\pi_\alpha \in \mathcal{D}$ defined in Equation (3.2) belong to $\mathcal{M}(m, \rho^2)$, and so they satisfy the convex log-Sobolev inequality. As a consequence the concentration inequality (3.38) holds for them.

Recall that (see Section 3.4.1 page 73) we need concentration bounds for the lower and upper tails of $\phi(Z)$, while the concentration inequality (3.38) does not contain these two sides.

Shu and Strzelecki (2017) gave a sufficient and necessary condition for a probability measure on the real line to satisfy the convex log-Sobolev inequality. They obtained concentration bounds for the lower and upper tails of convex functions of independent random variables which satisfy the convex log-Sobolev inequality.

The result obtained by Shu and Strzelecki (2017) allows us to state in Corollary 3.4.2 the appropriate concentration inequality for the probability distributions associated with the densities $\pi_\alpha \in \mathcal{D}$.

3.4.3.2 Concentration inequality for density π_α

In this Section we give the definition of the class of probability distributions $\mathcal{M}(m, \rho^2)$ and some of its properties. We show in Lemma 3.4.2 that the probability distributions associated with the densities $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)) belong to $\mathcal{M}(m, \rho^2)$, and so they satisfy the convex log-Sobolev inequality (3.37). Finally, we state in Corollary 3.4.2 the appropriate concentration inequality for our work which is a consequence of the concentration inequality stated in Corollary 1.7. of the paper by Shu and Strzelecki (2017).

Definition 3.4.2 (Definition 4 in Adameczak (2005)) For $m > 0$ and $\rho \geq 0$ let $\mathcal{M}(m, \rho^2)$ denote the class of probability distributions Π on \mathbb{R} for which

$$v^+(A) \leq \rho^2 \Pi(A),$$

for all sets A of the form $A = [x, \infty)$, $x \geq m$ and

$$v^-(A) \leq \rho^2 \Pi(A),$$

for all sets A of the form $A = (-\infty, -x]$, $x \geq m$, where v^+ is the measure on $[m, \infty)$ with density $x\Pi([x, \infty))$ and v^- is the measure on $(-\infty, -m]$ with density $-x\Pi((-\infty, x])$.

Example 3.4.1 (Example page 5 in Adamczak (2005)) The absolutely continuous distributions Π that satisfy for $t \geq m$,

$$\frac{d}{dt} \log \Pi([t, \infty)) \leq -\frac{t}{\rho^2} \quad \text{and} \quad \frac{d}{dt} \log \Pi((-\infty, -t]) \leq -\frac{t}{\rho^2}. \quad (3.39)$$

belong to $\mathcal{M}(m, \rho^2)$. In particular, if Π has density of the form $\exp(-V(x))$ with $dV(x)/dx \geq x/\rho^2$ and $dV(-x)/dx \leq -x/\rho^2$ then $\Pi \in \mathcal{M}(1, \rho^2)$.

It is shown by Adamczak (2005) that the probability distributions belonging to $\mathcal{M}(m, \rho^2)$ satisfy the convex log-Sobolev inequality (3.37). Let us denote by Π_α the probability distribution associated with the density $\pi_\alpha \in \mathcal{D}$ defined in Equation (3.2). In the following Lemma we will show that $\otimes \Pi_\alpha$ satisfies the convex log-Sobolev inequality (3.37).

Lemma 3.4.2 *There exists some m such that $\Pi_\alpha \in \mathcal{M}(m, \rho^2)$, and therefore $\otimes \Pi_\alpha$ satisfies the convex log-Sobolev inequality (3.37).*

The proof of Lemma 3.4.2 is given in Section 3.B.1 page 109.

As $\Pi_\alpha \in \mathcal{M}(m, \rho^2)$ and they satisfy the convex log-Sobolev inequality (3.37), so the concentration bound (3.38) holds for them. Recall that (see Section 3.4.1 page 73), we need a concentration bound for the both upper and lower tails of a convex function of the random variables that are distributed as Π_α . Therefore, the concentration bound (3.38) is not sufficient for our work. We state in Corollary 3.4.2 the appropriate concentration inequality for our work which is a consequence of the concentration inequality obtained by Shu and Strzelecki (2017). This result holds under a supplementary condition that we will state in the following Remark.

Remark 3.4.3 *Let Z be a random variable distributed as Π_α , then for every $s > 0$ the quantity $E(\exp(s|Z|))$ exists and is finite.*

The proof of Remark 3.4.3 is given in Section 3.B.2 page 110.

Note that, if $\alpha < 2$ then $E(\exp(s|Z|)) \not\leq \infty$.

Corollary 3.4.2 *Let $Z = (Z_1, \dots, Z_n)$ be i.i.d. random variables distributed as Π_α . Then there exists $A, B < \infty$ (depending only on C in the log-Sobolev inequality (3.37)), such that for any convex (or concave) function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ which is 1-Lipschitz (with respect to the Euclidean norm on \mathbb{R}^n) we have:*

$$P\left(|\phi(Z) - E(\phi(Z))| \geq u\right) \leq 2B \exp\left(-\frac{u^2}{8A}\right), \quad u \geq 0. \quad (3.40)$$

Corollary 3.4.2 is a consequence of the concentration inequality shown by Shu and Strzelecki (2017):

$$P\left(|\phi(Z) - M(\phi(Z))| \geq u\right) \leq B \exp\left(-\frac{u^2}{A}\right), \quad u \geq 0, \quad (3.41)$$

where M is the median of $\phi(Z)$.

The proof of Corollary 3.4.2 is given in Section 3.B.3 page 111 and is based on the fact that the concentration inequalities around the mean and the median are equivalent up to a numerical constant (Milman and Schechtman (1986)).

3.5 Proof of Theorem 3.3.1

The proof is based on four main lemmas proved in Section 3.5.2. In Section 3.5.1 other lemmas used all along the proof are stated.

Let us first establish inequalities that will be used in the following. Let $f \in \mathcal{H}$ and $v \in S_f$ (see (3.8)).

Using that for any $v \in S_f$, and any norm $\|\cdot\|$ in \mathcal{H}_v , $\|f_v\| - \|\widehat{f}_v\| \leq \|f_v - \widehat{f}_v\|$ and that for any $v \notin S_f$, $\|f_v\| = 0$, we get,

$$\sum_{v \in \mathcal{P}} \mu_v \|f_v\|_{\mathcal{H}_v} - \sum_{v \in \mathcal{P}} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} \leq \sum_{v \in S_f} \mu_v \|f_v - \widehat{f}_v\|_{\mathcal{H}_v} - \sum_{v \notin S_f} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v}, \quad (3.42)$$

and,

$$\sum_{v \in \mathcal{P}} \gamma_v \|f_v\|_n - \sum_{v \in \mathcal{P}} \gamma_v \|\widehat{f}_v\|_n \leq \sum_{v \in S_f} \gamma_v \|f_v - \widehat{f}_v\|_n - \sum_{v \notin S_f} \gamma_v \|\widehat{f}_v\|_n. \quad (3.43)$$

Combining (3.42), and (3.43), to the fact that for any function $f \in \mathcal{H}$, $\mathcal{L}(\widehat{f}) \leq \mathcal{L}(f)$, we obtain,

$$\|m - \widehat{f}\|_n^2 \leq \|m - f\|_n^2 + B,$$

with

$$B = 2V_{n,\varepsilon}(\widehat{f} - f) + \sum_{v \in S_f} [\mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v - f_v\|_n] - \sum_{v \notin S_f} [\mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v\|_n]. \quad (3.44)$$

If $\|m - f\|_n^2 \geq B$, we immediately get the result since in that case

$$\|m - \widehat{f}\|_n^2 \leq 2\|m - f\|_n^2 \leq 2\|m - f\|_n^2 + \sum_{v \in S_f} \mu_v + \sum_{v \in S_f} \gamma_v^2.$$

If $\|m - f\|_n^2 < B$, we get that

$$\|\widehat{f} - m\|_n^2 \leq 2B \quad (3.45)$$

$$\leq 4|V_{n,\varepsilon}(\widehat{f} - f)| + 2 \sum_{v \in S_f} [\mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v - f_v\|_n]. \quad (3.46)$$

The control of the empirical process $|V_{n,\varepsilon}(\widehat{f} - f)|$ is given by the following lemma (proved in Section 3.5.2.1, page 89).

Lemma 3.5.1 *Let $V_{n,\varepsilon}$ be defined in (3.21). For any f in \mathcal{F} , we consider the event \mathcal{T} defined as*

$$\mathcal{T} = \left\{ \forall f \in \mathcal{F}, \forall v \in \mathcal{P}, |V_{n,\varepsilon}(\widehat{f}_v - f_v)| \leq \kappa \lambda_{n,v}^2 \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \kappa \lambda_{n,v} \|\widehat{f}_v - f_v\|_n \right\}, \quad (3.47)$$

where $\lambda_{n,v}$ is defined in Equation (3.11) and where $\kappa = 10 + 4\Delta$. Then, for some positive constants c_1, c_2 ,

$$P_{X,\varepsilon}(\mathcal{T}) \geq 1 - c_1 \sum_{v \in \mathcal{P}} \exp(-nc_2 \lambda_{n,v}^2). \quad (3.48)$$

Conditioning on \mathcal{T} , Inequality (3.46) becomes

$$\begin{aligned} \|\widehat{f} - m\|_n^2 &\leq 4\kappa \sum_{v \in \mathcal{P}} [\lambda_{n,v}^2 \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \lambda_{n,v} \|\widehat{f}_v - f_v\|_n] + \\ &\quad 2 \sum_{v \in S_f} [\mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v - f_v\|_n], \end{aligned}$$

which may be decomposed as follows

$$\begin{aligned} \|\widehat{f} - m\|_n^2 &\leq \sum_{v \in S_f} [4\kappa \lambda_{n,v}^2 + 2\mu_v] \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \sum_{v \in S_f} [4\kappa \lambda_{n,v} + 2\gamma_v] \|\widehat{f}_v - f_v\|_n + \\ &\quad 4 \sum_{v \notin S_f} \kappa \lambda_{n,v}^2 \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + 4 \sum_{v \notin S_f} \kappa \lambda_{n,v} \|\widehat{f}_v - f_v\|_n. \end{aligned}$$

If we choose $C_1 \geq \kappa$ in Theorem 3.3.1, then $\kappa \lambda_{n,v}^2 \leq \mu_v$ and $\kappa \lambda_{n,v} \leq \gamma_v$ and the previous inequality becomes

$$\begin{aligned} \|\widehat{f} - m\|_n^2 &\leq 6 \sum_{v \in S_f} [\mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v - f_v\|_n] + \\ &\quad 4 \sum_{v \notin S_f} [\mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v\|_n]. \end{aligned} \quad (3.49)$$

Next we use the decomposability property of the penalty expressed in the following lemma (proved in Section 3.5.2.2 page 92).

Lemma 3.5.2 *For any $f \in \mathcal{F}$, under the assumptions of Theorem 3.3.1, conditionally on \mathcal{T} (see (3.47)), we have:*

$$\sum_{v \notin S_f} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} + \sum_{v \notin S_f} \gamma_v \|\widehat{f}_v\|_n \leq 3 \sum_{v \in S_f} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + 3 \sum_{v \in S_f} \gamma_v \|\widehat{f}_v - f_v\|_n. \quad (3.50)$$

Hence, by combining (3.49) and Lemma 3.5.2 we obtain

$$\|\widehat{f} - m\|_n^2 \leq 18 \sum_{v \in S_f} [\mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\widehat{f}_v - f_v\|_n].$$

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
86 in the regression model with non-Gaussian and non-bounded error

For each v , $\|\widehat{f}_v - f_v\|_{\mathcal{H}_v} \leq 2$ (because the functions \widehat{f}_v et f_v belong to the class \mathcal{F} , see (3.6)), and consequently, for some constant C ,

$$\|\widehat{f} - m\|_n^2 \leq C \left\{ \sum_{v \in S_f} \mu_v + \sum_{v \in S_f} \gamma_v \|\widehat{f}_v - f_v\|_n \right\}. \quad (3.51)$$

To finish the proof it remains to compare the two quantities $\sum_{v \in S_f} \|\widehat{f}_v - f_v\|_n^2$ and $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_n^2$. For that purpose we show that $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_n$ is less than $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_2^2$ plus an additive term coming from concentration results (see the Lemma given below). Next, thanks to the orthogonality of the spaces \mathcal{H}_v with respect to $L^2(P_X, \mathcal{X})$, $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_2^2 = \sum_{v \in S_f} \|\widehat{f}_v - f_v\|_2^2$. To conclude, it remains to consider several cases, according to the rankings of $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_2^2$ and $\|\sum_{v \in S_f} \widehat{f}_v - f_v\|_n^2$. This is the subject of the following lemma whose proof is given in Section 3.5.2.3, page 93.

Lemma 3.5.3 *For $f \in \mathcal{H}$, let \mathcal{A} be the event*

$$\mathcal{A} = \left\{ \forall f \in \mathcal{F}, \forall v \in \mathcal{P}, \|\widehat{f}_v - f_v\|_n \leq 2\|\widehat{f}_v - f_v\|_2 + \gamma_v \right\}. \quad (3.52)$$

Then, for some positive constant c_2 ,

$$P_{X,\varepsilon}(\mathcal{A}) \geq 1 - \sum_{v \in \mathcal{P}} \exp(-nc_2\gamma_v^2).$$

On the set \mathcal{A} , Inequality (3.51) provides that, for all $K > 0$

$$\begin{aligned} \frac{1}{C} \|\widehat{f} - m\|_n^2 &\leq \sum_{v \in S_f} [\mu_v + 2\gamma_v \|\widehat{f}_v - f_v\|_2 + \gamma_v^2], \\ &\leq \sum_{v \in S_f} [\mu_v + (1+K)\gamma_v^2 + \frac{1}{K} \|\widehat{f}_v - f_v\|_2^2], \end{aligned} \quad (3.53)$$

$$\begin{aligned} &\leq \sum_{v \in S_f} [\mu_v + (1+K)\gamma_v^2] + \frac{1}{K} \sum_{v \in \mathcal{P}} \|\widehat{f}_v - f_v\|_2^2, \\ &\leq \sum_{v \in S_f} [\mu_v + (1+K)\gamma_v^2] + \frac{1}{K} \left\| \sum_{v \in \mathcal{P}} \widehat{f}_v - f_v \right\|_2^2. \end{aligned} \quad (3.54)$$

Inequality (3.53) uses the inequality $2ab \leq \frac{1}{K}a^2 + Kb^2$ for all positive K , and Inequality (3.54) uses the orthogonality with respect to $L^2(P_X)$.

In the following we have to consider several cases, according to the rankings of $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2$ and $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_n$. More precisely, we consider two following cases:

Case 1: If $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2 \leq \|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_n$.

Case 2: If $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2 \geq \|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_n$.

Case 1: From (3.54), for any $f \in \mathcal{H}$, we get

$$\frac{1}{C} \|\widehat{f} - m\|_n^2 \leq \sum_{v \in S_f} [\mu_v + (1+K)\gamma_v^2] + \frac{1}{K} \|\widehat{f} - f\|_n^2.$$

Hence, using that for all $K' > 0$,

$$\|\widehat{f} - f\|_n^2 \leq (1+K') \|\widehat{f} - m\|_n^2 + (1 + \frac{1}{K'}) \|f - m\|_n^2, \quad (3.55)$$

we obtain for a suitable choice of K' , say $1+K' < K/C$, that, for some positive constant C' ,

$$\|\widehat{f} - m\|_n^2 \leq C' \left\{ \|f - m\|_n^2 + \sum_{v \in S_f} \mu_v + \sum_{v \in S_f} \gamma_v^2 \right\}.$$

This shows the result in Case 1.

Case 2: This case is solved by applying the following Lemma (proved in Section 3.5.2.4, page 93), which states that with high probability, $\|\widehat{f} - f\|_2 \leq \sqrt{2} \|\widehat{f} - f\|_n$.

Lemma 3.5.4 *Let $f = \sum_v f_v \in \mathcal{F}$ with support S_f , $\lambda_{n,v}$ be defined by (3.11), and let $\mathcal{G}(f)$ be the class of functions written as $g = \sum_{v \in \mathcal{P}} g_v$, such that $\|g_v\|_{\mathcal{H}_v} \leq 2$ satisfying for all $f \in \mathcal{F}$*

$$\begin{aligned} \mathbf{C1} \quad & \sum_{v \in \mathcal{P}} \mu_v \|g_v\|_{\mathcal{H}_v} + \sum_{v \in \mathcal{P}} \gamma_v \|g_v\|_n \leq 4 \sum_{v \in S_f} \mu_v \|g_v\|_{\mathcal{H}_v} + 4 \sum_{v \in S_f} \gamma_v \|g_v\|_n \\ \mathbf{C2} \quad & \sum_{v \in S_f} \gamma_v \|g_v\|_n \leq 2 \sum_{v \in S_f} \gamma_v \|g_v\|_2 + \sum_{v \in S_f} \gamma_v^2 \\ \mathbf{C3} \quad & \|g\|_n \leq \|g\|_2 \end{aligned}$$

Then the event

$$\left\{ \|g\|_n^2 \geq \frac{\|g\|_2^2}{2} \right\},$$

have probability greater than $1 - c_1 \exp(-nc_3 \sum_{v \in S_f} \lambda_{n,v}^2)$ for some constants c_1 and c_3 .

If f is such that $|S_f| = 0$, then Condition **C1** is not satisfied except if $g_v = 0$ for all $v \in \mathcal{P}$. Because we will apply Lemma 3.5.4 to $g_v = \widehat{f}_v - f_v$, this event has probability 0. If f is such that $|S_f| \geq 1$, then Condition **C1** is satisfied:

from Equation (3.50) in Lemma 3.5.2 we have,

$$\begin{aligned} & \sum_{v \notin S_f} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} + \sum_{v \in S_f} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \sum_{v \notin S_f} \gamma_v \|\widehat{f}_v\|_n + \sum_{v \in S_f} \gamma_v \|\widehat{f}_v - f_v\|_n \\ & \leq 3 \sum_{v \in S_f} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \sum_{v \in S_f} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} + 3 \sum_{v \in S_f} \gamma_v \|\widehat{f}_v - f_v\|_n + \sum_{v \in S_f} \gamma_v \|\widehat{f}_v\|_n, \\ & \Leftrightarrow \sum_{v \in \mathcal{P}} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + \sum_{v \in \mathcal{P}} \gamma_v \|\widehat{f}_v - f_v\|_n \leq 4 \sum_{v \in S_f} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + 4 \sum_{v \in S_f} \gamma_v \|\widehat{f}_v - f_v\|_n. \end{aligned}$$

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
88 in the regression model with non-Gaussian and non-bounded error

Moreover, Assumption $n\lambda_{n,v}^2 \geq -C_2 \log(\lambda_{n,v})$ implies that

$$\lambda_{n,v} = K_{n,v}/\sqrt{n} \text{ with } K_{n,v} \rightarrow \infty.$$

Then,

$$\exp(-nc_3 \sum_{v \in S_f} \lambda_{n,v}^2) \leq \exp(-c_3 |S_f| \min_{v \in \mathcal{P}} K_{n,v}^2),$$

and the event

$$\mathcal{C} = \left\{ \forall f \in \mathcal{F}, \text{ such that } g = \sum_{v \in \mathcal{P}} (\hat{f}_v - f_v) \in \mathcal{G}(f), \text{ and } \|g\|_n^2 \geq \frac{\|g\|_2^2}{2} \right\} \quad (3.56)$$

has probability greater than $1 - \eta/3$ for some $0 < \eta < 1$.

Conditioning on the events \mathcal{T} and \mathcal{A} (defined by (3.47) and (3.52)), $\sum_{v \in \mathcal{P}} (\hat{f}_v - f_v)$ belongs to the set $\mathcal{G}(f)$. According to (3.54), we conclude in the same way as in the first case.

Finally, it remains to quantify $P_{X,\varepsilon}(\mathcal{T} \cap \mathcal{A} \cap \mathcal{C})$. Following Lemma 3.5.1, and Lemma 3.5.3, \mathcal{T} , respectively \mathcal{A} , has probability greater than $1 - c_1 \sum_{v \in \mathcal{P}} \exp(-nc_2 \lambda_{n,v}^2)$, respectively $1 - \sum_{v \in \mathcal{P}} \exp(-n\gamma_v^2)$. Each of these probabilities is greater than $1 - \eta/3$ thanks to the assumption $n\lambda_{n,v}^2 \geq -C_2 \log \lambda_{n,v}$. □

3.5.1 Intermediate Lemmas

Lemma 3.5.5 *If $E_{X,\varepsilon}$ denotes the expectation with respect to the distribution of (X, ε) , we have for all $t > 0$,*

$$E_{X,\varepsilon} W_{n,2,v}(t) \leq Q_{n,v}(t).$$

Its proof is given in Section 3.5.3.1 page 101.

Lemma 3.5.6 *Let $b > 0$ and let $\mathcal{G}(t)$ be the following class of functions:*

$$\mathcal{G}(t) = \left\{ g_v \in \mathcal{H}_v, \|g_v\|_{\mathcal{H}_v} \leq 2, \|g_v\|_2 \leq t, \|g_v\|_\infty \leq b \right\}. \quad (3.57)$$

Let $\Omega_{v,t}$ be the event defined as

$$\Omega_{v,t} = \left\{ \sup_{g_v \in \mathcal{G}(t)} \{ \|g_v\|_2 - \|g_v\|_n \} \leq \frac{bt}{2} \right\}. \quad (3.58)$$

Then for any $t \geq \nu_{n,v}$, the event $\Omega_{v,t}$ has probability greater than $1 - \exp(-c_2 nt^2)$, for some positive constant c_2 .

Its proof is given in Section 3.5.3.2, page 101.

Lemma 3.5.7 *For any function $g_v \in \mathcal{H}_v$ satisfying $\|g_v\|_{\mathcal{H}_v} \leq 2$, $\|g_v\|_\infty \leq b$ and $\|g_v\|_2 \geq t$, for all $t \geq \nu_{n,v}$ and $b \geq 1$, the event*

$$\left(1 - \frac{b}{2}\right) \|g_v\|_2 \leq \|g_v\|_n \leq \left(1 + \frac{b}{2}\right) \|g_v\|_2$$

has probability greater than $1 - \exp(-c_2 nt^2)$ for some positive constant c_2 .

Its proof is given in Section 3.5.3.3, page 103.

Lemma 3.5.8 *If E_ε denotes the expectation with respect to the distribution of ε , we have*

$$P_{X,\varepsilon}\left(|W_{n,n,v}(t) - E_\varepsilon(W_{n,n,v}(t))| \geq \delta t\right) \leq 2B \exp\left(-\frac{n\delta^2}{8A}\right). \quad (3.59)$$

Its proof is given in Section 3.5.3.4, page 103.

Lemma 3.5.9 *Conditionally on the space $\Omega_{v,t}$ defined by (3.58), we have the following inequalities:*

$$P_{X,\varepsilon}\left(|W_{n,2,v}(t) - E_\varepsilon(W_{n,2,v}(t))| \geq \delta t\right) \leq 2B \exp\left(-\frac{n\delta^2}{32A}\right), \quad (3.60)$$

$$P_X\left(E_\varepsilon W_{n,2,v}(t) - E_{X,\varepsilon}(W_{n,2,v}(t)) \geq x\right) \leq \exp\left(-\frac{nx^2}{Q_{n,v}(t)}\right). \quad (3.61)$$

Its proof is given in Section 3.5.3.5, page 104.

Lemma 3.5.10 *Let $\lambda_{n,v}$ be defined at Equation (3.11), Δ at Equation (3.10) and $\kappa = 10 + 4\Delta$. Conditionally on the space $\Omega_{v,\lambda_{n,v}}$ defined at Equation (3.58), for some positive constants c_1, c_2 , with probability greater than $1 - c_1 \exp(-c_2 n \lambda_{n,v}^2)$, we have*

$$W_{n,n,v}(\lambda_{n,v}) \leq \kappa \lambda_{n,v}^2 \quad \text{and} \quad E_\varepsilon W_{n,n,v}(\lambda_{n,v}) \leq \kappa \lambda_{n,v}^2. \quad (3.62)$$

Its proof is given in Section 3.5.3.6, page 106.

3.5.2 Proof of lemma 3.5.1 to 3.5.4

3.5.2.1 Proof of lemma 3.5.1

For $f \in \mathcal{F}$ and $v \in \mathcal{P}$, let $g_v = \widehat{f}_v - f_v$. Note that $\|g_v\|_{\mathcal{H}_v} \leq 2$. Let us show that

$$|V_{n,\varepsilon}(g_v)| \leq \kappa \left(\lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + \lambda_{n,v} \|g_v\|_n \right). \quad (3.63)$$

We start by writing that

$$|V_{n,\varepsilon}(g_v)| = \|g_v\|_{\mathcal{H}_v} \left| V_{n,\varepsilon}\left(\frac{g_v}{\|g_v\|_{\mathcal{H}_v}}\right) \right| \leq \|g_v\|_{\mathcal{H}_v} W_{n,n,v}\left(\frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}}\right). \quad (3.64)$$

Consider the two following cases:

Case A: $\|g_v\|_n \leq \lambda_{n,v} \|g_v\|_{\mathcal{H}_v}$,

Case B: $\|g_v\|_n > \lambda_{n,v} \|g_v\|_{\mathcal{H}_v}$.

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
90 in the regression model with non-Gaussian and non-bounded error**

Case A: Since $\|g_v\|_n \leq \lambda_{n,v} \|g_v\|_{\mathcal{H}_v}$, we have

$$W_{n,n,v} \left(\frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}} \right) \leq W_{n,n,v}(\lambda_{n,v}).$$

We then apply Lemma 3.5.10, page 89, and conclude that (3.63) holds in Case A for each $v \in \mathcal{P}$ since, with high probability

$$|V_{n,\varepsilon}(g_v)| \leq \kappa \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} \leq \kappa \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + \kappa \lambda_{n,v} \|g_v\|_n. \quad (3.65)$$

Case B: Consider now the case $\|g_v\|_n > \lambda_{n,v} \|g_v\|_{\mathcal{H}_v}$ and let us show that for any $v \in \mathcal{P}$,

$$W_{n,n,v} \left(\frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}} \right) \leq \kappa \lambda_{n,v} \|g_v\|_n.$$

Let r_v be a deterministic number such that $r_v > \lambda_{n,v}$. Our first step relies on the study of the process $W_{n,n,v}(r_v)$, for $r_v > \lambda_{n,v}$. In that case we state two results:

R1 For any deterministic $r_v \geq \lambda_{n,v}$, with probability greater than $1 - c_1 \exp(-c_2 n \lambda_{n,v}^2)$,

$$W_{n,n,v}(r_v) \leq \kappa r_v \lambda_{n,v}. \quad (3.66)$$

R2 Inequality (3.66) continues to hold for random r_v of the form

$$r_v = \frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}}.$$

Combining these two points implies that, with probability greater than $1 - c_1 \exp(-c_2 n \lambda_{n,v}^2)$,

$$\|g_v\|_{\mathcal{H}_v} W_{n,n,v} \left(\frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}} \right) \leq \kappa \|g_v\|_n \lambda_{n,v}.$$

Consequently, in Case B, according to (3.64), for each v , Inequality (3.63) holds because

$$|V_{n,\varepsilon}(g_v)| \leq \kappa \|g_v\|_n \lambda_{n,v} \leq \kappa \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + \kappa \lambda_{n,v} \|g_v\|_n.$$

This ends up the proof of Lemma 3.5.1.

Proof of R1 From Lemma 3.5.8, page 89 with $t = r_v$ and $\delta = \lambda_{n,v}$, we get that with probability greater than $1 - 2B \exp(-n \lambda_{n,v}^2 / 8A)$,

$$W_{n,n,v}(r_v) \leq E_\varepsilon(W_{n,n,v}(r_v)) + r_v \lambda_{n,v} \quad (3.67)$$

Next we prove that for some positive r_v , with probability greater than $1 - \exp(-nc \lambda_{n,v}^2)$, we have

$$E_\varepsilon(W_{n,n,v}(r_v)) \leq \kappa r_v \lambda_{n,v}. \quad (3.68)$$

Let $\widehat{\nu}_{n,v}$ be defined as the smallest solution of $E_\varepsilon(W_{n,n,v}(t)) \leq \kappa t^2$. For $W_{n,n,v}$, defined by (3.23), we write

$$E_\varepsilon(W_{n,n,v}(r_v)) = \frac{r_v}{\widehat{\nu}_{n,v}} E_\varepsilon \sup \left\{ |V_{n,\varepsilon}(g_v)|, \|g_v\|_{\mathcal{H}_v} \leq 2\left(\frac{\widehat{\nu}_{n,v}}{r_v}\right), \|g_v\|_n \leq \widehat{\nu}_{n,v} \right\}.$$

Besides, Lemma 3.5.10 stated that on the event $\Omega_{v,\lambda_{n,v}}$, $E_\varepsilon(W_{n,n,v}(\lambda_{n,v})) \leq \kappa \lambda_{n,v}^2$. It follows from the definition of $\widehat{\nu}_{n,v}$, and Lemma 3.5.6, that $\widehat{\nu}_{n,v} \leq \lambda_{n,v}$ for all $v \in \mathcal{P}$ with probability greater than $1 - \exp(-nc_2 \sum_{v \in \mathcal{P}} \lambda_{n,v}^2)$. Consequently, for any deterministic r_v such that $r_v \geq \lambda_{n,v}$, we have

$$\widehat{\nu}_{n,v} \leq \lambda_{n,v} \leq r_v \Leftrightarrow \frac{\widehat{\nu}_{n,v}}{r_v} \leq 1,$$

and so,

$$\begin{aligned} E_\varepsilon(W_{n,n,v}(r_v)) &= \frac{r_v}{\widehat{\nu}_{n,v}} E_\varepsilon \sup \left\{ |V_{n,\varepsilon}(g_v)|, \|g_v\|_{\mathcal{H}_v} \leq 2, \|g_v\|_n \leq \widehat{\nu}_{n,v} \right\}, \\ &\leq \frac{r_v}{\widehat{\nu}_{n,v}} E_\varepsilon(W_{n,n,v}(\widehat{\nu}_{n,v})) \leq \frac{r_v}{\widehat{\nu}_{n,v}} \kappa \widehat{\nu}_{n,v}^2 = \kappa r_v \widehat{\nu}_{n,v} \leq \kappa r_v \lambda_{n,v}. \end{aligned}$$

Proof of R2 Let us prove **R2** by using a peeling-type argument. Our aim is to prove that (3.66) holds for any r_v of the form

$$r_v = \frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}}.$$

Since $\|g_v\|_\infty / \|g_v\|_{\mathcal{H}_v} \leq 1$, we have $\|g_v\|_n / \|g_v\|_{\mathcal{H}_v} \leq 1$. We thus restrict ourselves to r_v satisfying $r_v = \|g_v\|_n / \|g_v\|_{\mathcal{H}_v}$ with $\|g_v\|_n / \|g_v\|_{\mathcal{H}_v} \in (\lambda_{n,v}, 1]$.

We start by splitting the interval $(\lambda_{n,v}, 1]$ into M disjoint intervals such that

$$(\lambda_{n,v}, 1] = \cup_{k=1}^M (2^{k-1} \lambda_{n,v}, 2^k \lambda_{n,v}],$$

for some M that will be chosen later. Consider the event \mathcal{D}^c defined as follows:

$$\mathcal{D}^c = \left\{ \exists v \in \mathcal{P} \text{ and } \exists \bar{g}_v, \text{ such that } |V_{n,\varepsilon}(\bar{g}_v)| \geq \kappa \lambda_{n,v} \|\bar{g}_v\|_n, \text{ with } \frac{\|\bar{g}_v\|_n}{\|\bar{g}_v\|_{\mathcal{H}_v}} \in (\lambda_{n,v}, 1] \right\}.$$

We prove that, for some positive constants c_1, c_2 ,

$$P(\mathcal{D}^c) \leq c_1 \exp(-c_2 n \lambda_{n,v}^2).$$

For $\bar{g}_v \in \mathcal{D}^c$, let \bar{k} be the integer in $\{1, \dots, M\}$, such that

$$2^{\bar{k}-1} \lambda_{n,v} \leq \frac{\|\bar{g}_v\|_n}{\|\bar{g}_v\|_{\mathcal{H}_v}} \leq 2^{\bar{k}} \lambda_{n,v}.$$

This \bar{k} satisfies

$$\|\bar{g}_v\|_{\mathcal{H}_v} W_{n,n,v}(2^{\bar{k}} \lambda_{n,v}) \geq \|\bar{g}_v\|_{\mathcal{H}_v} W_{n,n,v}\left(\frac{\|g_v\|_n}{\|g_v\|_{\mathcal{H}_v}}\right) \geq |V_{n,\varepsilon}(\bar{g}_v)| \geq \kappa \lambda_{n,v} \|\bar{g}_v\|_n.$$

Therefore, we get

$$W_{n,n,v}(2^{\bar{k}}\lambda_{n,v}) \geq \kappa\lambda_{n,v} \frac{\|\bar{g}_v\|_n}{\|\bar{g}_v\|_{\mathcal{H}_v}} \geq \kappa\lambda_{n,v}^2 2^{\bar{k}-1} \geq \kappa \frac{\lambda_{n,v}}{2} 2^{\bar{k}}\lambda_{n,v}.$$

By taking $r_v = 2^{\bar{k}}\lambda_{n,v}$ in (3.66), we have

$$\mathcal{P}\left(W_{n,n,v}(2^{\bar{k}}\lambda_{n,v}) \geq \kappa \frac{\lambda_{n,v}}{2} 2^{\bar{k}}\lambda_{n,v}\right) \leq c_1 \exp(-c_2 n \lambda_{n,v}^2).$$

Now let us write \mathcal{D}^c as follows:

$$\mathcal{D}^c = \bigcup_{k=1}^M \left\{ \exists v \text{ and } \exists \bar{g}_v \text{ such that } |V_{n,\varepsilon}(\bar{g}_v)| \geq \kappa\lambda_{n,v}\|\bar{g}_v\|_n, \text{ with } \frac{\|\bar{g}_v\|_n}{\|\bar{g}_v\|_{\mathcal{H}}} \in (2^{k-1}\lambda_{n,v}, 2^k\lambda_{n,v}] \right\}.$$

The set \mathcal{D}^c has probability smaller than $c_1 M \exp(-c_2 n \lambda_{n,v}^2)$. If we choose M such that $\log M \leq (c_2/2)n\lambda_{n,v}^2$, then the probability of the set \mathcal{T} is greater than

$$1 - \sum_{v \in \mathcal{P}} c_1 \exp\left(-\frac{c_2}{2} n \lambda_{n,v}^2\right).$$

It follows that **R2** is proved which ends up the proof of Lemma 3.5.1. \square

3.5.2.2 Proof of lemma 3.5.2

Starting from (3.45) with B defined by Equation (3.44), we write

$$\begin{aligned} \frac{1}{2}\|\hat{f} - m\|_n^2 &\leq 2|V_{n,\varepsilon}(\hat{f} - f)| + \sum_{v \in S_f} [\mu_v \|\hat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\hat{f}_v - f_v\|_n] - \\ &\quad \sum_{v \notin S_f} [\mu_v \|\hat{f}_v\|_{\mathcal{H}_v} + \gamma_v \|\hat{f}_v\|_n]. \end{aligned}$$

On the event \mathcal{T} defined in (3.47) we have

$$\begin{aligned} \frac{1}{2}\|\hat{f} - m\|_n^2 &\leq 2\kappa \sum_{v \in \mathcal{P}} \lambda_{n,v}^2 \|\hat{f}_v - f_v\|_{\mathcal{H}_v} + 2\kappa \sum_{v \in \mathcal{P}} \lambda_{n,v} \|\hat{f}_v - f_v\|_n + \\ &\quad \sum_{v \in S_f} [\mu_v \|\hat{f}_v - f_v\|_{\mathcal{H}_v} + \gamma_v \|\hat{f}_v - f_v\|_n] - \sum_{v \notin S_f} [\mu_v \|\hat{f}_v\|_{\mathcal{H}_v} + \gamma_v \|\hat{f}_v\|_n]. \end{aligned}$$

Rearranging the terms we obtain that

$$\begin{aligned} \frac{1}{2}\|\hat{f} - m\|_n^2 &\leq \sum_{v \in S_f} (2\kappa\lambda_{n,v}^2 + \mu_v) \|\hat{f}_v - f_v\|_{\mathcal{H}_v} + \sum_{v \in S_f} (2\kappa\lambda_{n,v} + \gamma_v) \|\hat{f}_v - f_v\|_n + \\ &\quad \sum_{v \notin S_f} (2\kappa\lambda_{n,v}^2 - \mu_v) \|\hat{f}_v\|_{\mathcal{H}_v} + \sum_{v \notin S_f} (2\kappa\lambda_{n,v} - \gamma_v) \|\hat{f}_v\|_n. \end{aligned}$$

Now, thanks to Assumption (3.12) with $C_1 \geq \kappa$ we have $\kappa\lambda_{n,v}^2 \leq \mu_v$ and $2\kappa\lambda_{n,v} \leq \gamma_v$ and Lemma 3.5.2 is shown since

$$0 \leq \frac{1}{2} \|\widehat{f} - m\|_n^2 \leq 3 \sum_{v \in S_f} \mu_v \|\widehat{f}_v - f_v\|_{\mathcal{H}_v} + 3 \sum_{v \in S_f} \|\widehat{f}_v - f_v\|_n - \sum_{v \notin S_f} \mu_v \|\widehat{f}_v\|_{\mathcal{H}_v} - \sum_{v \notin S_f} \gamma_v \|\widehat{f}_v\|_n.$$

□

3.5.2.3 Proof of lemma 3.5.3

Let us consider the following two cases:

- ✓ $\|\widehat{f}_v - f_v\|_2 \leq \gamma_v$. We apply Lemma 3.5.6 (page 88) to the function $g_v = \widehat{f}_v - f_v$. It satisfies $g_v \in \mathcal{G}(\gamma_v)$ with $b = 2$ (recall that $\|\cdot\|_\infty \leq \|\cdot\|_{\mathcal{H}_v}$). Moreover, $\gamma_v \geq C_1\lambda_{n,v} \geq C_1\nu_{nv} \geq \nu_{n,v}$ as soon as $C_1 \geq 1$.

It follows that, for some positive c_2 , with probability greater than $1 - \exp(-nc_2\gamma_v^2)$,

$$\|\widehat{f}_v - f_v\|_n \leq \|\widehat{f}_v - f_v\|_2 + \gamma_v.$$

- ✓ $\|\widehat{f}_v - f_v\|_2 \geq \gamma_v$. We apply Lemma 3.5.7 (page 88) to the function $g_v = \widehat{f}_v - f_v$ with $b = 2$. It follows that, for some positive c_2 , with probability greater than $1 - \exp(-nc_2\gamma_v^2)$,

$$\|\widehat{f}_v - f_v\|_n \leq 2\|\widehat{f}_v - f_v\|_2.$$

□

3.5.2.4 Proof of lemma 3.5.4

Throughout the proof, we make use of the quantity d_n defined as follows:

For $\beta < 1/\alpha$ and some constant η' ,

$$d_n^2 \geq \eta' n^{\alpha\beta-1}. \quad (3.69)$$

Let $\mathcal{G}(f)$ and $\mathcal{G}'(f)$ be the following sets:

$$\mathcal{G}(f) = \left\{ g = \sum_{v \in \mathcal{P}} g_v, \text{ satisfying } \|g_v\|_{\mathcal{H}_v} \leq 2, \text{ and Conditions } \mathbf{C1}, \mathbf{C2}, \mathbf{C3} \right\},$$

$$\mathcal{G}'(f) = \left\{ g \in \mathcal{G}(f), \text{ such that } \|g\|_2 = d_n \right\}.$$

In order to prove this lemma we consider two cases: if $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2 \geq d_n$, and if $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2 \leq d_n$.

First, we suppose that $\|\sum_{v \in \mathcal{P}} \widehat{f}_v - f_v\|_2 \geq d_n$, and we consider the two events \mathcal{B} and \mathcal{B}' defined as follows:

$$\mathcal{B} = \left\{ \forall h \in \mathcal{G}, \|h\|_n^2 \geq \frac{\|h\|_2^2}{2}, \text{ and } \|h\|_2 \geq d_n \right\},$$

and

$$\mathcal{B}' = \left\{ \forall h \in \mathcal{G}', \|h\|_n^2 \geq \frac{d_n^2}{2} \right\}. \quad (3.70)$$

If $h \in \mathcal{B}'$, then $h \in \mathcal{G}$, $\|h\|_2 = d_n$ and $\|h\|_n^2 \geq d_n^2/2$. It follows that $\|h\|_n^2 \geq \|h\|_2^2/2$ and $\|h\|_2 \geq d_n$. We just showed that the event \mathcal{B}' is included into the event \mathcal{B} . So, this case is proved if the event \mathcal{B}' holds with high probability. Consider

$$Z_n(\mathcal{G}') = \sup_{g \in \mathcal{G}'} \left\{ d_n^2 - \|g\|_n^2 \right\}.$$

We show that the event $Z_n(\mathcal{G}') \leq d_n^2/2$ has probability greater than $1 - c_1 \exp(-nc_3 d_n^2)$.

Consider a $d_n/8$ -covering of $(\mathcal{G}', \|\cdot\|_n)$. So that, for all g in \mathcal{G}' there exists g^k such that

$$\|g - g^k\|_n \leq \frac{d_n}{8}.$$

The associated proper covering number is:

$$N_{\text{pr}} = N_{\text{pr}}\left(\frac{d_n}{8}, \mathcal{G}', \|\cdot\|_n\right). \quad (3.71)$$

Now, for all $g \in \mathcal{G}'$, we write:

$$d_n^2 - \|g\|_n^2 = T_1 + T_2, \quad (3.72)$$

with $T_1 = \|g^k\|_n^2 - \|g\|_n^2$ and $T_2 = d_n^2 - \|g^k\|_n^2$. The proof is splitted into four steps:

Step 1 The first step consists in showing that

$$T_1 = \|g^k\|_n^2 - \|g\|_n^2 \leq \frac{d_n^2}{4}. \quad (3.73)$$

Step 2 The second step consists in proving that, for N_{pr} given at Equation (3.71) and for some constant C ,

$$P_X \left(\max_{k \in \{1, \dots, N_{\text{pr}}\}} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4} \right) \leq \exp \left(\log N_{\text{pr}} - Cnd_n^2 \right).$$

Step 3 The third step concerns the control of N_{pr} . Let σ_α^2 be the variance of a random variable distributed with density $\pi_\alpha \in \mathcal{D}$ (see Equation (3.2)), then for some $K > 0$,

$$\begin{aligned} \frac{1}{K} \log N_{\text{pr}} \leq & \left(32\sigma_\alpha \sqrt{n} (E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|) / d_n \right)^\alpha \mathbf{1}_{(0, 32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|)}(d_n) + \\ & \mathbf{1}_{[32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)}(d_n). \end{aligned}$$

Step 4 The last step consists in bounding from above the Gaussian complexity. For some $\kappa > 0$

$$E_\varepsilon \sup_{g \in \mathcal{G}'} \sum_{v \in \mathcal{P}} |V_{n,\varepsilon}(g_v)| \leq \frac{4\kappa}{C_1} \left\{ \sum_{v \in S_f} (2\mu_v + \gamma_v^2) + 2 \left(\sum_{v \in S_f} \gamma_v^2 \right)^{\frac{1}{2}} d_n \right\},$$

Let us conclude the proof of the lemma before proving these four steps.

Putting together Steps 3 and 4 we have:

If $d_n \in [32\sigma_\alpha\sqrt{n}E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)$, then

$$\frac{1}{K} \log N_{\text{pr}}\left(\frac{d_n}{8}, \mathcal{G}', \|\cdot\|_n\right) \leq 1.$$

Thanks to Step 2,

$$P_X\left(T_2 \geq \frac{d_n^2}{4}\right) \leq P_X\left(\max_{k \in \{1, \dots, N_{\text{pr}}\}} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4}\right) \leq K \exp\left(-Cnd_n^2\right),$$

and, therefore

$$P_X\left(Z_n(\mathcal{G}') \leq \frac{d_n^2}{2}\right) \leq K \exp\left(-Cnd_n^2\right). \quad (3.74)$$

If $d_n \in (0, 32\sigma_\alpha\sqrt{n}E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|]$, then

$$\begin{aligned} \frac{1}{K} \log N_{\text{pr}}\left(\frac{d_n}{8}, \mathcal{G}', \|\cdot\|_n\right) &\leq (32\sigma_\alpha)^\alpha n^{\frac{\alpha}{2}} \left(\frac{E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{d_n}\right)^\alpha, \\ &\leq (32\sigma_\alpha)^\alpha n^{\frac{\alpha}{2}} \left(\frac{4\kappa}{C_1 d_n} \left(\sum_{v \in S_f} (2\mu_v + \gamma_v^2) + 2\left(\sum_{v \in S_f} \gamma_v^2\right)^{\frac{1}{2}} d_n\right)\right)^\alpha, \\ &\leq \left(\frac{128\kappa\sigma_\alpha}{C_1}\right)^\alpha n^{\frac{\alpha}{2}} \left(\frac{\sum_{v \in S_f} (2\mu_v + \gamma_v^2)}{d_n} + 2\left(\sum_{v \in S_f} \gamma_v^2\right)^{\frac{1}{2}}\right)^\alpha. \end{aligned}$$

We have to show that $\log N_{\text{pr}} - Cnd_n^2 \leq -c_3nd_n^2$ or equivalently that $\log N_{\text{pr}} \leq \tilde{C}nd_n^2$, where $\tilde{C} = C - c_3$.

Let $A = K(128\kappa\sigma_\alpha/C_1)^\alpha$. We have,

$$\begin{aligned} \log N_{\text{pr}} \leq \tilde{C}nd_n^2 &\Leftrightarrow An^{\frac{\alpha}{2}} \left(\frac{\sum_{v \in S_f} (2\mu_v + \gamma_v^2)}{d_n} + 2\left(\sum_{v \in S_f} \gamma_v^2\right)^{\frac{1}{2}}\right)^\alpha \leq \tilde{C}nd_n^2, \\ &\Leftrightarrow \left(\frac{\sum_{v \in S_f} (2\mu_v + \gamma_v^2)}{d_n} + 2\left(\sum_{v \in S_f} \gamma_v^2\right)^{\frac{1}{2}}\right)^\alpha \leq \frac{\tilde{C}}{A} n^{1-\frac{\alpha}{2}} d_n^2, \\ &\Leftrightarrow \frac{\sum_{v \in S_f} (2\mu_v + \gamma_v^2)}{d_n} + 2\left(\sum_{v \in S_f} \gamma_v^2\right)^{\frac{1}{2}} \leq \left(\frac{\tilde{C}}{A}\right)^{\frac{1}{\alpha}} n^{\frac{1}{\alpha}-\frac{1}{2}} d_n^{\frac{2}{\alpha}}. \end{aligned}$$

Because $\gamma_v = C_1\lambda_{n,v}$ and $\mu_v = C_1\lambda_{n,v}^2$,

$$\log N_{\text{pr}} \leq \tilde{C}nd_n^2 \Leftrightarrow C_1(2 + C_1) \frac{\sum_{v \in S_f} \lambda_{n,v}^2}{d_n} + 2C_1 \left(\sum_{v \in S_f} \lambda_{n,v}^2\right)^{\frac{1}{2}} \leq \left(\frac{\tilde{C}}{A}\right)^{\frac{1}{\alpha}} n^{\frac{1}{\alpha}-\frac{1}{2}} d_n^{\frac{2}{\alpha}}.$$

Considering the first term in the left hand side, let

$$B = \frac{1}{2} \times \frac{1}{C_1(2 + C_1)} \left(\frac{\tilde{C}}{A}\right)^{\frac{1}{\alpha}},$$

Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
96 in the regression model with non-Gaussian and non-bounded error

then

$$\frac{\sum_{v \in S_f} \lambda_{n,v}^2}{d_n} \leq B n^{\frac{1}{\alpha} - \frac{1}{2}} d_n^{\frac{2}{\alpha}} \Leftrightarrow d_n^2 \geq B^{-\frac{2\alpha}{2+\alpha}} \left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{2\alpha}{2+\alpha}} n^{\frac{\alpha-2}{\alpha+2}}.$$

As $\sum_{v \in S_f} \lambda_{n,v}^2 \leq C_3 n^{2\beta-1}$ (see Equation (3.14)), we get

$$B^{-\frac{2\alpha}{2+\alpha}} \left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{2\alpha}{2+\alpha}} n^{\frac{\alpha-2}{\alpha+2}} \leq \left(\frac{B}{C_3} \right)^{-\frac{2\alpha}{2+\alpha}} n^{\frac{4\alpha\beta}{2+\alpha} - 1}.$$

Therefore, the inequality

$$C_1(2 + C_1) \frac{\sum_{v \in S_f} \lambda_{n,v}^2}{d_n} \leq \frac{1}{2} \left(\frac{\tilde{C}}{A} \right)^{\frac{1}{\alpha}} n^{\frac{1}{\alpha} - \frac{1}{2}} d_n^{\frac{2}{\alpha}},$$

will be satisfied if

$$d_n^2 \geq \left(\frac{C_3}{B} \right)^{\frac{2\alpha}{\alpha+2}} n^{\frac{4\alpha\beta}{\alpha+2} - 1}.$$

For the second term, let

$$B' = \frac{1}{2} \times \frac{1}{2C_1} \left(\frac{\tilde{C}}{A} \right)^{\frac{1}{\alpha}},$$

then

$$\left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{1}{2}} \leq B' n^{\frac{1}{\alpha} - \frac{1}{2}} d_n^{\frac{2}{\alpha}} \Leftrightarrow d_n^2 \geq B'^{-\alpha} \left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{\alpha}{2}} n^{\frac{\alpha-2}{2}}.$$

As $\sum_{v \in S_f} \lambda_{n,v}^2 \leq C_3 n^{2\beta-1}$ (see Equation (3.14)), then

$$B'^{-\alpha} \left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{\alpha}{2}} n^{\frac{\alpha-2}{2}} \leq \left(\frac{C_3}{B'^2} \right)^{\frac{\alpha}{2}} n^{\alpha\beta-1}.$$

Therefore the inequality

$$2C_1 \left(\sum_{v \in S_f} \lambda_{n,v}^2 \right)^{\frac{1}{2}} \leq \frac{1}{2} \left(\frac{\tilde{C}}{A} \right)^{\frac{1}{\alpha}} n^{\frac{1}{\alpha} - \frac{1}{2}} d_n^{\frac{2}{\alpha}},$$

will be satisfied if

$$d_n^2 \geq \left(\frac{C_3}{B'^2} \right)^{\frac{\alpha}{2}} n^{\alpha\beta-1}.$$

As $\alpha > 2$, $4\alpha\beta/(\alpha + 2) < \alpha\beta$. Therefore, there exists a constant η' , take for example

$$\eta' = \max \left(\left(\frac{C_3}{B'^2} \right)^{\frac{\alpha}{2}}, \left(\frac{C_3}{B} \right)^{\frac{2\alpha}{\alpha+2}} \right),$$

such that if $d_n^2 \geq \eta' n^{\alpha\beta-1}$, then $\log N_{\text{pr}} \leq \tilde{C} n d_n^2$, and Step 2 states that

$$P_X \left(T_2 \geq \frac{d_n^2}{4} \right) \leq P_X \left(\max_{k \in \{1, \dots, N_{\text{pr}}\}} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4} \right) \leq \exp \left(-c_3 n d_n^2 \right).$$

Now, we have

$$P_X \left(Z_n(\mathcal{G}') \leq \frac{d_n^2}{2} \right) = P_X \left(\max_{g^1, \dots, g^N} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4} \right) \leq \exp \left(-c_3 n d_n^2 \right). \quad (3.75)$$

Finally, we obtain for $c_1 = \max(K, 1)$ and $c_3 \leq C$ (see Equations (3.74) and (3.75)):

$$P_X \left(Z_n(\mathcal{G}') \leq \frac{d_n^2}{2} \right) \leq c_1 \exp \left(-c_3 n d_n^2 \right).$$

Moreover, for n large enough, we have $\sum_{v \in S_f} \lambda_{n,v}^2 \leq d_n^2 \leq 1$ (see Equations (3.14) and (3.69)), and

$$1 - c_1 \exp \left(-c_3 n d_n^2 \right) \geq 1 - c_1 \exp \left(-c_3 n \sum_{v \in S_f} \lambda_{n,v}^2 \right).$$

Therefore,

$$P_X \left(Z_n(\mathcal{G}') \leq \frac{d_n^2}{2} \right) \leq c_1 \exp \left(-c_3 n \sum_{v \in S_f} \lambda_{n,v}^2 \right).$$

Before proving the Steps 1 to 4 let us solve the second case: if $\|\sum_{v \in \mathcal{P}} \hat{f}_v - f_v\|_2 \leq d_n$ then we consider the event \mathcal{B}'' defined as follows:

$$\mathcal{B}'' = \left\{ \forall h \in \mathcal{G}, \|h\|_n^2 \geq \frac{\|h\|_2^2}{2}, \text{ and } \|h\|_2 \leq d_n \right\}.$$

We have that the event \mathcal{B}' defined in Equation (3.70) is included in \mathcal{B}'' and the same proof as in the first case applies.

Proofs of Steps 1 to 4 The proofs of Step 1 and Step 2 are strictly the same as in the Gaussian case. More precisely

Proof of Step 1:

It is easy to see that,

$$\begin{aligned} T_1 &= \|g^k\|_n^2 - \|g\|_n^2 = \frac{1}{n} \sum_{i=1}^n [(g^k(X_i))^2 - (g(X_i))^2] \\ &= \frac{1}{n} \sum_{i=1}^n [g^k(X_i) - g(X_i)][g^k(X_i) + g(X_i)] \\ &\leq \|g^k - g\|_n \left(\frac{1}{n} \sum_{i=1}^n [g^k(X_i) + g(X_i)]^2 \right)^{\frac{1}{2}} \end{aligned}$$

where in the inequality above we used Cauchy Schwarz inequality. Using the inequality $(a + b)^2 \leq 2a^2 + 2b^2$, $g \in \mathcal{G}'$, and the property that g satisfies Condition **C3**, we get

$$\frac{1}{n} \sum_{i=1}^n [g^k(X_i) + g(X_i)]^2 \leq 2\|g^k\|_n^2 + 2\|g\|_n^2 \leq 4d_n^2.$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
98 in the regression model with non-Gaussian and non-bounded error**

Besides, the covering set is constructed such that $\|g^k - g\|_n \leq d_n/8$. It follows that Step 1 is proved.

Proof of Step 2:

We prove that for some constant C ,

$$P_X\left(T_2 \geq \frac{d_n^2}{4}\right) \leq P_X\left(\max_{1 \leq k \leq N_{\text{pr}}} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4}\right) \leq \exp(\log N_{\text{pr}} - Cnd_n^2).$$

As $g^k \in \mathcal{G}'$, $d_n = \|g^k\|_2$. Then

$$\max_{1 \leq k \leq N_{\text{pr}}} [d_n^2 - \|g^k\|_n^2] = \max_{1 \leq k \leq N_{\text{pr}}} [\|g^k\|_2^2 - \|g^k\|_n^2].$$

Applying Theorem 3.5. in [Chung and Lu \(2006\)](#) with $X = \sum_i (g^k(X_i))^2$, for all positive λ we have:

$$P_X\left(\sum_{i=1}^n [g^k(X_i)]^2 \leq n\mathbb{E}(g^k(X_i))^2 - \lambda\right) \leq \exp\left(-\frac{\lambda^2}{2n\mathbb{E}(g^k(X))^4}\right),$$

or equivalently,

$$P_X\left(\|g^k\|_2^2 - \|g^k\|_n^2 \geq \frac{\lambda}{n}\right) \leq \exp\left(-\frac{\lambda^2}{2n\mathbb{E}(g^k(X))^4}\right).$$

Taking $\lambda = nd_n^2/4$ and using that $\|g^k\|_2^2 = d_n^2$ we get

$$P_X\left(d_n^2 - \|g^k\|_n^2 \geq \frac{d_n^2}{4}\right) \leq \exp\left(-\frac{nd_n^4}{32\mathbb{E}(g^k(X))^4}\right).$$

It follows that

$$\begin{aligned} P_X\left(\max_{1 \leq k \leq N_{\text{pr}}} [d_n^2 - \|g^k\|_n^2] \geq \frac{d_n^2}{4}\right) &\leq \sum_{k=1}^{N_{\text{pr}}} \exp\left(-\frac{nd_n^4}{32\mathbb{E}(g^k(X))^4}\right) \\ &\leq \exp\left(\log N_{\text{pr}} - \frac{nd_n^4}{32 \max_k \mathbb{E}(g^k(X))^4}\right). \end{aligned} \quad (3.76)$$

Moreover, $g \in \mathcal{H}$, so $g = \sum_{v \in \mathcal{P}} g_v$, where the functions g_v are centered and orthogonal in $L^2(P_X)$. Therefore $\mathbb{E}(g(X))^4$ is the sum of the following terms:

$$\begin{aligned} A_1 &= \sum_{v \in \mathcal{P}} E_X g_v^4(X_v), \\ A_2 &= \binom{4}{2} \sum_{v \neq v'} E_X g_v^2(X_v) g_{v'}^2(X_{v'}), \\ A_3 &= \binom{4}{3} \sum_{v_1 \neq v_2 \neq v_3} E_X g_{v_1}^2(X_{v_1}) g_{v_2}(X_{v_2}) g_{v_3}(X_{v_3}), \\ A_4 &= \binom{4}{3} \sum_{v_1 \neq v_2} E_X g_{v_1}^3(X_{v_1}) g_{v_2}(X_{v_2}), \\ A_5 &= \binom{4}{1} \sum_{v_1 \neq v_2 \neq v_3 \neq v_4} E_X g_{v_1}(X_{v_1}) g_{v_2}(X_{v_2}) g_{v_3}(X_{v_3}) g_{v_4}(X_{v_4}). \end{aligned}$$

Using the Cauchy Schwartz inequality and the fact that $\|g_v\|_\infty \leq \|g_v\|_{\mathcal{H}_v} \leq 2$, and $\|g\|_2 = d_n$ (because $g \in \mathcal{G}'$), we get that A_1 is proportional to d_n^2 , A_2, A_3, A_5 to d_n^4 , and A_4 to d_n^3 . For example,

$$A_1 = \sum_{v \in \mathcal{P}} E_X g_v^4(X_v) \leq \|g\|_\infty^2 \sum_{v \in \mathcal{P}} \|g_v\|_2^2 = \|g\|_\infty^2 \sum_{v \in \mathcal{P}} g_v^2 \leq 4d_n^2.$$

After calculation of the terms A_i , since d_n^2 is assumed to be smaller than one, we get that:

$$\max_k E_X (g^k(X))^4 \leq C d_n^2 (1 + O(d_n^2)). \quad (3.77)$$

Step 2 is proved by combining (3.76) and (3.77).

We now focus on Step 3 and Step 4:

Proof of Step 3:

Let N_{pr} be defined at Equation (3.71). We prove that

$$\begin{aligned} \frac{1}{K} \log N_{\text{pr}} \left(\frac{d_n}{8}, \mathcal{G}', \|\cdot\|_n \right) &\leq \left(32\sigma_\alpha \sqrt{n} (E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|) / d_n \right)^\alpha \mathbf{1}_{(0, 32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|]}(d_n) + \\ &\quad \mathbf{1}_{[32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)}(d_n). \end{aligned}$$

We start from Equation (3.27) and write that:

$$\log N_{\text{pr}} \left(\frac{d_n}{8}, \mathcal{G}', \|\cdot\|_n \right) \leq \log N \left(\frac{d_n}{16}, \mathcal{G}', \|\cdot\|_n \right).$$

Next, we use Corollary 3.4.1:

Let $Z = (Z_1, \dots, Z_n)$ be i.i.d. random variables distributed with density $\pi_\alpha \in \mathcal{D}$ defined in Equation (3.2) with $\sqrt{\text{var}(Z_i)} = \sigma_\alpha$. Set $T = \mathcal{G}'$, $\delta = \sqrt{n}d_n/16$ and $M = n \times E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|$, then for all $\alpha \geq 2$ we have,

$$\begin{aligned} \log N \left(\frac{d_n}{16}, \mathcal{G}', \|\cdot\|_n \right) &= \log N \left(\frac{\sqrt{n}d_n}{16}, \mathcal{G}', \|\cdot\| \right), \\ &\leq K \left(\frac{32n E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|}{\sqrt{n}d_n} \right)^\alpha \mathbf{1}_{(0, 2n \times E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|]} \left(\frac{\sqrt{n}d_n}{16} \right) + \\ &\quad K \left(\frac{32n E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|}{\sqrt{n}d_n} \right)^2 \mathbf{1}_{[2n \times E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|, \infty)} \left(\frac{\sqrt{n}d_n}{16} \right), \end{aligned}$$

or equivalently,

$$\begin{aligned} \log N \left(\frac{d_n}{16}, \mathcal{G}', \|\cdot\|_n \right) &\leq K \left(\frac{32n E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|}{\sqrt{n}d_n} \right)^\alpha \mathbf{1}_{(0, 32\sqrt{n} E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|]}(d_n) + \\ &\quad K \left(\frac{32n E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|}{\sqrt{n}d_n} \right)^2 \mathbf{1}_{[32\sqrt{n} E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)|, \infty)}(d_n). \end{aligned}$$

Take $\varepsilon_i = Z_i/\sigma_\alpha = h(Z_i)$ for $i = 1, \dots, n$, then $\text{var}(\varepsilon_i) = 1$ and,

$$E_\varepsilon(\varepsilon_i) = E_\varepsilon(h(Z_i)) = \int h(Z_i) \pi_\alpha(Z_i) dZ_i = \frac{1}{\sigma_\alpha} \int Z_i \pi_\alpha(Z_i) dZ_i = \frac{1}{\sigma_\alpha} E_Z(Z_i).$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
100 in the regression model with non-Gaussian and non-bounded error**

Therefore, $E_Z \sup_{g \in \mathcal{G}'} |V_{n,Z}(g)| = \sigma_\alpha E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|$ and,

$$\begin{aligned} \log N\left(\frac{d_n}{16}, \mathcal{G}', \|\cdot\|_n\right) &\leq K \left(\frac{32n\sigma_\alpha E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{\sqrt{nd_n}}\right)^\alpha \mathbf{1}_{(0, 32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|]}(d_n) + \\ &\quad K \left(\frac{32n\sigma_\alpha E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{\sqrt{nd_n}}\right)^2 \mathbf{1}_{[32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)}(d_n), \\ &\leq K (32\sigma_\alpha)^\alpha n^{\frac{\alpha}{2}} \left(\frac{E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{d_n}\right)^\alpha \mathbf{1}_{(0, 32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|]}(d_n) + \\ &\quad K \left(\frac{32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{d_n}\right)^2 \mathbf{1}_{[32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)}(d_n), \\ &\leq K (32\sigma_\alpha)^\alpha n^{\frac{\alpha}{2}} \left(\frac{E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|}{d_n}\right)^\alpha \mathbf{1}_{(0, 32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|]}(d_n) + \\ &\quad K \mathbf{1}_{[32\sigma_\alpha \sqrt{n} E_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|, \infty)}(d_n). \end{aligned}$$

Proof of Step 4:

This Step consists in bounding from above the quantity $\mathbb{E}_\varepsilon \sup_{g \in \mathcal{G}'} |V_{n,\varepsilon}(g)|$. According to Inequality (3.63) we have,

$$\sum_{v \in \mathcal{P}} |V_{n,\varepsilon}(g_v)| \leq \kappa \left\{ \sum_{v \in \mathcal{P}} \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + \sum_{v \in \mathcal{P}} \lambda_{n,v} \|g_v\|_n \right\},$$

with $\lambda_{n,v}$ defined by Equation (3.11) satisfying Equation (3.12) for all $v \in \mathcal{P}$. It follows

$$\begin{aligned} \sup_{g \in \mathcal{G}'} \sum_{v \in \mathcal{P}} |V_{n,\varepsilon}(g_v)| &\leq \kappa \sup_{g \in \mathcal{G}'} \left\{ \sum_{v \in \mathcal{P}} \lambda_{n,v}^2 \|g_v\|_{\mathcal{H}_v} + \sum_{v \in \mathcal{P}} \lambda_{n,v} \|g_v\|_n \right\}, \\ &\leq \frac{\kappa}{C_1} \sup_{g \in \mathcal{G}'} \left\{ \sum_{v \in \mathcal{P}} \mu_v \|g_v\|_{\mathcal{H}_v} + \sum_{v \in \mathcal{P}} \gamma_v \|g_v\|_n \right\}. \end{aligned}$$

Thanks to Condition **C1** and using $\|g_v\|_{\mathcal{H}_v} \leq 2$ we obtain then:

$$\begin{aligned} \sup_{g \in \mathcal{G}'} \sum_{v \in \mathcal{P}} |V_{n,\varepsilon}(g_v)| &\leq \frac{4\kappa}{C_1} \left\{ \sup_{g \in \mathcal{G}'} \sum_{v \in S_f} \mu_v \|g_v\|_{\mathcal{H}_v} + \sup_{g \in \mathcal{G}'} \sum_{v \in S_f} \gamma_v \|g_v\|_n \right\}, \\ &\leq \frac{4\kappa}{C_1} \left\{ 2 \sum_{v \in S_f} \mu_v + \sup_{g \in \mathcal{G}'} \sum_{v \in S_f} \gamma_v \|g_v\|_n \right\}. \end{aligned}$$

Now, according to Condition **C2**, we get

$$\begin{aligned} \sup_{g \in \mathcal{G}'} \sum_{v \in \mathcal{P}} |V_{n,\varepsilon}(g_v)| &\leq \frac{4\kappa}{C_1} \left\{ 2 \sum_{v \in S_f} \mu_v + 2 \sup_{g \in \mathcal{G}'} \sum_{v \in S_f} \gamma_v \|g_v\|_2 + \sum_{v \in S_f} \gamma_v^2 \right\}, \\ &\leq \frac{4\kappa}{C_1} \left\{ \sum_{v \in S_f} (2\mu_v + \gamma_v^2) + 2 \sup_{g \in \mathcal{G}'} \left(\sum_{v \in S_f} \gamma_v^2 \right)^{1/2} \left(\sum_{v \in S_f} \|g_v\|_2^2 \right)^{1/2} \right\}, \\ &\leq \frac{4\kappa}{C_1} \left\{ \sum_{v \in S_f} (2\mu_v + \gamma_v^2) + 2 \left(\sum_{v \in S_f} \gamma_v^2 \right)^{1/2} d_n \right\}, \end{aligned}$$

where in the second inequality we used Cauchy Schwarz inequality and the third inequality coming from the fact that for all $g \in \mathcal{G}'$, $\|g\|_2^2 = d_n^2 \geq \sum_{v \in S_f} \|g_v\|_2^2$. \square

3.5.3 Proofs of intermediate Lemmas

3.5.3.1 Proof of Lemma 3.5.5

The kernel k_v is written as :

$$k_v(X_v, X'_v) = \sum_{\ell \geq 1} \omega_{v,\ell} \phi_{v,\ell}(X_v) \phi_{v,\ell}(X'_v)$$

where $\{\phi_{v,\ell}\}_{\ell=1}^{\infty}$ is an orthonormal basis of $L^2(P_v)$ with $P_v = \prod_{a \in v} P_a$.

Let us consider the class of functions $\mathcal{K}(t)$ defined as

$$\mathcal{K}(t) = \{g_v \in \mathcal{H}_v, \|g_v\|_{\mathcal{H}_v} \leq 2, \|g_v\|_2 \leq t\}.$$

It comes that

$$g_v = \sum_{\ell} a_{\ell} \phi_{v,\ell}, \quad \text{with } \|g_v\|_{\mathcal{H}_v}^2 = \sum_{\ell} \frac{a_{\ell}^2}{\omega_{v,\ell}} \leq 4, \quad \text{and } \|g_v\|_2^2 = \sum_{\ell} a_{\ell}^2 \leq t^2$$

In the following, we set $\mu_{v,\ell}(t) = \min\{t^2, \omega_{v,\ell}\}$. Hence

$$\sum_{\ell} \frac{a_{\ell}^2}{\mu_{v,\ell}(t)} \leq \frac{1}{t^2} \sum_{\ell} a_{\ell}^2 + \sum_{\ell} \frac{a_{\ell}^2}{\omega_{v,\ell}} = \frac{1}{t^2} \|g_v\|_2^2 + \|g_v\|_{\mathcal{H}_v}^2 \leq 5, \quad (3.78)$$

as soon as $g_v \in \mathcal{K}(t)$.

Now, let us prove the lemma:

$$\begin{aligned} E_{X,\varepsilon} W_{n,2,v}(t) &= E_{X,\varepsilon} \sup_{g \in \mathcal{K}(t)} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \sum_{\ell} a_{\ell} \phi_{v,\ell}(X_{vi}) \right|, \\ &= E_{X,\varepsilon} \sup_{g \in \mathcal{K}(t)} \left| \frac{1}{n} \sum_{\ell} \frac{a_{\ell}}{\sqrt{\mu_{v,\ell}(t)}} \sum_{i=1}^n \varepsilon_i \sqrt{\mu_{v,\ell}(t)} \phi_{v,\ell}(X_{vi}) \right|, \\ &\leq \sqrt{5} \sqrt{E_{X,\varepsilon} \sum_{\ell} \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i \sqrt{\mu_{v,\ell}(t)} \phi_{v,\ell}(X_{vi}) \right)^2}. \end{aligned}$$

The last inequality follows from the Cauchy-Schwartz inequality and Inequality (3.78).

Now, simple calculation leads to

$$E_{X,\varepsilon} W_{n,2,v}(t) \leq \sqrt{5} \sqrt{\frac{1}{n} \sum_{\ell} \mu_{v,\ell}(t)}.$$

□

3.5.3.2 Proof of Lemma 3.5.6

Using that $|\sqrt{a} - \sqrt{b}| \leq \sqrt{|a - b|}$, we get

$$|\|g_v\|_2 - \|g_v\|_n| \leq \sqrt{|\|g_v\|_2^2 - \|g_v\|_n^2|}.$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
102 in the regression model with non-Gaussian and non-bounded error**

Hence

$$\left\{ \|g_v\|_\infty \leq b, \left| \|g_v\|_2 - \|g_v\|_n \right| \geq \frac{bt}{2} \right\} \subset \left\{ \left| \|g_v\|_2^2 - \|g_v\|_n^2 \right| \geq \frac{b^2 t^2}{4} \right\}.$$

The centered process

$$\left| \|g_v\|_2^2 - \|g_v\|_n^2 \right| = \left| \frac{1}{n} \sum_{i=1}^n g_v^2(X_{v,i}) - \mathbb{E}(g_v^2(X_v)) \right|,$$

satisfies a concentration inequality given, for example, by Theorem 2.1 in [Bartlett et al. \(2005\)](#) : if \mathcal{C} is a class of functions f such that $\|f\|_\infty \leq B$ and $Ef(X) = 0$, and if there exists $\gamma > 0$ such that for every $f \in \mathcal{C}$, $\text{Var}f(X) \leq \gamma^2$. Then for every $x > 0$, with probability at least $1 - e^{-x}$,

$$\sup_{f \in \mathcal{C}} \frac{1}{n} \left| \sum_{j=1}^n f(X_j) \right| \leq \inf_{\alpha > 0} \left\{ 2(1 + \alpha) E \left(\sup_{f \in \mathcal{C}} \frac{1}{n} \left| \sum_{j=1}^n f(X_j) \right| \right) + \sqrt{\frac{2x}{n}} \gamma + B \left(\frac{1}{3} + \frac{1}{\alpha} \right) \frac{x}{n} \right\}. \quad (3.79)$$

For any $t > 0$, for $\mathcal{G}(t)$ defined by (3.57), let us consider the class of functions $\mathcal{C}(t)$ defined as follows

$$\mathcal{C}(t) = \left\{ f \text{ such that } f = g_v^2 - \mathbb{E}(g_v^2), \text{ with } g_v \in \mathcal{G}(t) \right\}.$$

Note that if $f \in \mathcal{C}(t)$, $E_X f(X_v) = 0$ and $\|f\|_\infty \leq b^2$. We have to study

$$\gamma^2(t) = \sup_{g_v \in \mathcal{G}(t)} E_X \left(g_v^2(X) - \|g_v\|_2^2 \right)^2 \text{ and } \Gamma(t) = E_X \left(\sup_{g_v \in \mathcal{G}(t)} \left| \|g_v\|_n^2 - \|g_v\|_2^2 \right| \right).$$

It is easy to see that

$$\gamma^2(t) \leq b^2 \sup_{g_v \in \mathcal{G}(t)} E_X (g_v(X) + \|g_v\|_2)^2 \leq 4b^2 t^2.$$

Let ζ_i be i.i.d. Rademacher random variables and let $E_{X,\zeta}$ denotes the expectation with respect to the law of (X, ζ) . By a symmetrization argument,

$$\Gamma(t) \leq 2E_{X,\zeta} \sup_{g_v \in \mathcal{G}(t)} \left| \frac{1}{n} \sum_{i=1}^n \zeta_i g_v^2(X_i) \right|.$$

Since $\|g_v\|_\infty \leq b$, applying the contraction principal (see [Ledoux and Talagrand \(1991\)](#)) we get that, for $Q_{n,v}(t)$ defined by (3.9),

$$E_{X,\zeta} \sup_{g_v \in \mathcal{G}(t)} \left| \frac{1}{n} \sum_{i=1}^n \zeta_i g_v^2(X_i) \right| \leq 4b E_{X,\zeta} \sup_{g_v \in \mathcal{G}(t)} \left| \frac{1}{n} \sum_{i=1}^n \zeta_i g_v(X_i) \right| \leq 4b Q_{n,v}(t).$$

The last inequality was proved by [Mendelson \(2002\)](#), Theorem 41 (see the proof of Lemma 3.5.5). Now, thanks to (3.79) we get that for all $x > 0$, with probability greater than $1 - e^{-x}$

$$\sup_{g_v \in \mathcal{G}(t)} \left| \|g_v\|_n^2 - \|g_v\|_2^2 \right| \leq \inf_{\alpha > 0} \left\{ 16(1 + \alpha)b Q_{n,v}(t) + \sqrt{\frac{2x}{n}} 2bt + b^2 \left(\frac{1}{3} + \frac{1}{\alpha} \right) \frac{x}{n} \right\}.$$

Taking $x = c_2 n t^2$, $t \geq \nu_n$, we have that with probability greater than $1 - e^{-c_2 n t^2}$

$$\sup_{g_v \in \mathcal{G}(t)} \left| \|g_v\|_n^2 - \|g_v\|_2^2 \right| \leq \inf_{\alpha > 0} t^2 \left\{ 16(1 + \alpha)b\Delta + \sqrt{2c_2}4b + b^2 \left(\frac{1}{3} + \frac{1}{\alpha} \right) c_2 \right\}.$$

The infimum of the right hand side is reached in $\alpha = \sqrt{c_2 b / 16\Delta}$, and equals

$$\frac{b^2 c_2}{3} + 8\sqrt{\Delta c_2} b^{3/2} + 4(4\Delta + \sqrt{2c_2})b.$$

The constants Δ and c_2 should satisfy that this infimum is strictly smaller than $b^2/4$. For example, if $16\Delta < b/8$, it remains to choose c_2 small enough such that

$$b \left(\frac{c_2}{3} + \frac{\sqrt{2c_2}}{2} \right) + 4\sqrt{2c_2} < \frac{b}{8}.$$

□

3.5.3.3 Proof of Lemma 3.5.7

Let $t > \nu_{n,v}$ and h be defined as

$$h = \frac{t g_v}{\|g_v\|_2}.$$

If g_v satisfies the assumptions of the lemma, then h satisfies $\|h\|_2 = t$, $\|h\|_{\mathcal{H}} \leq 2$ and $\|h\|_{\infty} \leq b$. Applying Lemma 3.5.6 (page 88) to the function h , we obtain that for all $t \geq \nu_{n,v}$, with probability greater than $1 - \exp(-c_2 n t^2)$, we have

$$|t - \|h\|_n| \leq \frac{bt}{2} \quad \text{for all } h \in \mathcal{G}(t).$$

This concludes the proof of the lemma.

□

3.5.3.4 Proof of Lemma 3.5.8

We apply Corollary 3.4.2 to

$$\phi(\varepsilon_1, \dots, \varepsilon_n) = \frac{\sqrt{n}}{t} W_{n,n,v}(t).$$

Using Cauchy-Schwarz Inequality and the fact that $\|g_v\|_n \leq t$,

$$|\phi(\varepsilon) - \phi(\varepsilon')| \leq \frac{\sqrt{n}}{t} \sup_{\|g_v\|_n \leq t} \|g_v\|_n \|\varepsilon - \varepsilon'\|_n \leq \frac{\sqrt{n}}{t} t \|\varepsilon - \varepsilon'\|_n,$$

leading to $\|\phi\|_L = 1$. So,

$$P_{X,\varepsilon} \left(\left| \frac{\sqrt{n}}{t} W_{n,n,v}(t) - \frac{\sqrt{n}}{t} E_{\varepsilon} W_{n,n,v}(t) \right| \geq u \right) \leq 2B \exp \left(-\frac{u^2}{8A} \right),$$

and Lemma 3.5.8 is proved by taking $\delta = u/\sqrt{n}$.

□

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
104 in the regression model with non-Gaussian and non-bounded error**

3.5.3.5 Proof of Lemma 3.5.9

We start with the proof of (3.60) in Lemma 3.5.9 by applying once again Corollary 3.4.2, to the function

$$\phi(\varepsilon) = \phi(\varepsilon_1, \dots, \varepsilon_n) = \frac{\sqrt{n}}{2t} W_{n,2,v}(t).$$

On the event $\Omega_{v,t}$ defined by (3.58), we have

$$\|g_v\|_n \leq \frac{bt}{2} + \|g_v\|_2.$$

Besides if $\|g_v\|_{\mathcal{H}_v} \leq 2$, then $\|g_v\|_\infty \leq 2$. Therefore applying Lemma 3.5.6 with $b = 2$, we get that if $\|g_v\|_2 \leq t$,

$$|\phi(\varepsilon) - \phi(\varepsilon')| \leq \frac{\sqrt{n}}{2t} \sup_{\|g_v\|_n \leq 2t} \|g_v\|_n \|\varepsilon - \varepsilon'\|_n \leq \frac{\sqrt{n}}{2t} 2t \|\varepsilon - \varepsilon'\|_n,$$

leading to $\|\phi\|_L = 1$. So,

$$P_{X,\varepsilon} \left(\left\{ \left| \frac{\sqrt{n}}{2t} W_{n,2,v}(t) - \frac{\sqrt{n}}{2t} E_\varepsilon(W_{n,2,v}(t)) \right| \geq u \right\} \cap \Omega_{v,t}^c \right) \leq 2B \exp \left(-\frac{u^2}{8A} \right),$$

and inequality (3.60) in Lemma 3.5.9 is proved by taking $\delta = 2u/\sqrt{n}$.

We now come to the proof of the inequality (3.61) in Lemma 3.5.9 using a Poissonian inequality for self-bounded processes (see Boucheron et al. (2000)) and Theorem 5.6, p 158 in Massart and Picard (2007)). Let us recall it in the particular case we are interested in:

Theorem 3.5.1 *Let X_1, \dots, X_n be n i.i.d. random variables. For $i \in \{1, \dots, n\}$ let*

$$X_{(-i)} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n).$$

Let h be a non-negative and bounded measurable function of $X = (X_1, \dots, X_n)$. Assume that for all $i \in \{1, \dots, n\}$, there exists a measurable function h_i of $X_{(-i)}$ such that $0 < h - h_i \leq 1$, and $\sum_{i=1}^n (h - h_i) \leq h$. Then, for all $x > 0$, we have

$$P \left(h \geq E(h) + x \right) \leq \exp \left(-\frac{x^2}{2E(h)} \right).$$

We apply this result to h defined as

$$h = h(X_1, \dots, X_n) = nE_\varepsilon W_{n,2,v}(t) = nE_\varepsilon \sup \left\{ |V_{n,\varepsilon}(g_v)|, \|g_v\|_2 \leq t, \|g_v\|_{\mathcal{H}_v} \leq 2 \right\}.$$

The variable h is positive, and because the distribution of $(\varepsilon_1, \dots, \varepsilon_n)$ is symmetric, we have that

$$h = E_\varepsilon \sup \left\{ nV_{n,\varepsilon}(g_v), \|g_v\|_2 \leq t, \|g_v\|_{\mathcal{H}_v} \leq 2 \right\}.$$

Let τ be the function in \mathcal{H}_v such that $h = E_\varepsilon n V_{n,\varepsilon}(\tau)$ (note that τ depends on (X_1, \dots, X_n) and on $(\varepsilon_1, \dots, \varepsilon_n)$), and let

$$h_i = E_\varepsilon \sup_{g_v} \sum_{j \neq i} \varepsilon_j g_v(X_j).$$

We show that h and h_i satisfy the assumptions of Theorem 3.5.1:

$$\begin{aligned} h - h_i &= E_\varepsilon \left(\varepsilon_i \tau(X_i) + \sum_{j \neq i} \varepsilon_j \tau(X_j) - \sup_{g_v} \sum_{j \neq i} \varepsilon_j g_v(X_j) \right), \\ &\leq E_\varepsilon \left(\varepsilon_i \tau(X_i) \right), \\ &\leq E_\varepsilon \left(|\varepsilon_i| \sup_{x \in \mathcal{X}} |\tau(X)| \right), \\ &\leq 2E_\varepsilon \left(|\varepsilon_i| \right), \end{aligned}$$

where the last inequality comes from the fact that $\sup_{x \in \mathcal{X}} |\tau(X)| \leq \|\tau\|_{\mathcal{H}_v} \leq 2$.

Let $Z = (Z_1, \dots, Z_n)$ be i.i.d. random variables distributed with density $\pi_\alpha \in \mathcal{D}$ defined in Equation (3.2) with $\sqrt{\text{var}(Z_i)} = \sigma_\alpha$. Take $\varepsilon_i = Z_i/\sigma_\alpha$ for $i = 1, \dots, n$, then $\text{var}(\varepsilon_i) = 1$ and $E_\varepsilon(|\varepsilon_i|) = E_Z(|Z_i|)/\sigma_\alpha$. We have:

$$E_Z(|Z_i|) = \int_{\mathbb{R}} |Z_i| a_\alpha \exp(-|Z_i|^\alpha) dZ_i.$$

Take $|Z_i| = u^{1/\alpha}$ and

$$dZ_i = \begin{cases} \frac{1}{\alpha} u^{\frac{1}{\alpha}-1} du & \text{if } Z_i \geq 0, \\ -\frac{1}{\alpha} u^{\frac{1}{\alpha}-1} du & \text{if } Z_i \leq 0. \end{cases}$$

Therefore,

$$\begin{aligned} E_Z(|Z_i|) &= \int_0^{+\infty} a_\alpha u^{\frac{1}{\alpha}} \exp(-u) \frac{1}{\alpha} u^{\frac{1}{\alpha}-1} du - \int_{+\infty}^0 a_\alpha u^{\frac{1}{\alpha}} \exp(-u) \frac{1}{\alpha} u^{\frac{1}{\alpha}-1} du, \\ &= 2 \int_0^{+\infty} a_\alpha u^{\frac{1}{\alpha}} \exp(-u) \frac{1}{\alpha} u^{\frac{1}{\alpha}-1} du, \\ &= a_\alpha \int_0^{+\infty} \frac{2}{\alpha} u^{\frac{2}{\alpha}-1} \exp(-u) du, \\ &= a_\alpha \frac{2}{\alpha} \Gamma\left(\frac{2}{\alpha}\right) = a_\alpha \Gamma\left(1 + \frac{2}{\alpha}\right), \end{aligned}$$

where $\Gamma(\cdot)$ is the gamma function.

It follows that,

$$h - h_i \leq \frac{2a_\alpha}{\sigma_\alpha} \Gamma\left(1 + \frac{2}{\alpha}\right).$$

Moreover, $h - h_i \geq 0$ since

$$h = E_\varepsilon \left(\sup_{g_v} \sum_{j=1}^n \varepsilon_j g_v(X_j) \right) = E_\varepsilon \left(E_{\varepsilon_i} \sup_{g_v} \sum_{j=1}^n \varepsilon_j g_v(X_j) \right) \geq E_\varepsilon \left(\sup_{g_v} E_{\varepsilon_i} \sum_{j=1}^n \varepsilon_j g_v(X_j) \right) = h_i.$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
106 in the regression model with non-Gaussian and non-bounded error**

Finally we have:

$$\sum_i (h - h_i) = \sum_{i=1}^n E_\varepsilon \left(\varepsilon_i \tau(X_i) + \sum_{j \neq i}^n \varepsilon_j \tau(X_j) - \sup_{g-v} \sum_{j \neq i}^n \varepsilon_j g_v(X_j) \right) \leq \sum_{i=1}^n E_\varepsilon \varepsilon_i \tau(X_i) = h.$$

Therefore, following Theorem 3.5.1, we get that for all positive u

$$P_{X,\varepsilon} \left(E_\varepsilon W_{n,2,v}(t) - E_{X,\varepsilon} W_{n,2,v}(t) \leq \frac{u}{n} \right) \leq \exp \left(- \frac{u^2}{E_{X,\varepsilon} W_{n,2,v}(t)} \right).$$

As $E_{X,\varepsilon} W_{n,2,v}(t) \leq Q_{n,v}(t)$, see Lemma 3.5.5 page 88, we get the expected result since for all positive x

$$P_X \left(E_\varepsilon W_{n,2,v}(t) \geq E_{X,\varepsilon} W_{n,2,v}(t) + x \right) \leq \exp \left(- \frac{nx^2}{Q_{n,v}(t)} \right).$$

□

3.5.3.6 Proof of Lemma 3.5.10

From Lemma 3.5.8, page 89 with $t = \lambda_{n,v} = \delta$, with probability greater than $1 - 2B \exp(-n\lambda_{n,v}^2/8A)$, we get that:

$$E_\varepsilon(W_{n,n,v}(\lambda_{n,v})) \leq \lambda_{n,v}^2 + W_{n,n,v}(\lambda_{n,v}) \quad (3.80)$$

The next step consists in comparing $W_{n,n,v}(\lambda_{n,v})$ and $W_{n,2,v}(2\lambda_{n,v})$. Recall that $\lambda_{n,v} \geq \nu_{n,v}$, see (3.11). Let g_v such that $\|g_v\|_n \leq \lambda_{n,v}$.

- ✓ When $\|g_v\|_2 \leq \lambda_{n,v}$, according to Lemma 3.5.6 (page 88), taking $b = 2$, since since $\|g_v\|_n \leq \lambda_{n,v}$, we get that with probability greater than $1 - \exp(-c_2 n \lambda_{n,v}^2)$,

$$\|g_v\|_n - \lambda_{n,v} \leq \|g_v\|_2 \leq \|g_v\|_n + \lambda_{n,v} \leq 2\lambda_{n,v}.$$

- ✓ When $\|g_v\|_2 \geq t$, we apply Lemma 3.5.7 (page 88) with $b = 2$. For any function g_v such that $\|g_v\|_\infty \leq 2$, and $\|g_v\|_2 \geq \lambda_{n,v}$, we have $\|g_v\|_2 \leq 2\|g_v\|_n \leq 2\lambda_{n,v}$.

This implies that, with probability greater than $1 - \exp(-c_2 n \lambda_{n,v}^2)$ we have

$$W_{n,n,v}(\lambda_{n,v}) \leq W_{n,2,v}(2\lambda_{n,v}).$$

We now study the process $W_{n,2,v}(\lambda_{n,v})$. By applying (3.60) in Lemma 3.5.9, page 89, with $\delta = t = \lambda_{n,v}$ we get that with probability greater than $1 - 2B \exp(-n\lambda_{n,v}^2/32A)$

$$W_{n,2,v}(\lambda_{n,v}) \leq \lambda_{n,v}^2 + E_\varepsilon(W_{n,2,v}(\lambda_{n,v})).$$

It follows that

$$\begin{aligned} E_\varepsilon W_{n,n,v}(\lambda_{n,v}) &\leq \lambda_{n,v}^2 + W_{n,n,v}(\lambda_{n,v}), \\ &\leq \lambda_{n,v}^2 + W_{n,2,v}(2\lambda_{n,v}), \\ &\leq 5\lambda_{n,v}^2 + E_\varepsilon(W_{n,2,v}(2\lambda_{n,v})). \end{aligned}$$

Next, we apply (3.61) in Lemma 3.5.9, with $t = 2\lambda_{n,v}$ and $x = 4\lambda_{n,v}^2$. We get that

$$E_\varepsilon W_{n,2,v}(2\lambda_{n,v}) \leq 4\lambda_{n,v}^2 + E_{X,\varepsilon}(W_{n,2,v}(2\lambda_{n,v})),$$

with probability greater than

$$1 - 2 \exp\left(-16 \frac{n\lambda_{n,v}^4}{Q_{n,v}(2\lambda_{n,v})}\right) \geq 1 - 2 \exp\left(-\frac{4n\lambda_{n,v}^2}{\Delta}\right).$$

The last inequality comes from the definition of $\nu_{n,v}$, see (3.10), and from the fact that $\lambda_{n,v} \geq \nu_{n,v}$, see (3.11).

Putting everything together, we get that with probability greater than $1 - c_1 \exp(-c_2 n \lambda_{n,v}^2)$ for some positive constants c_1, c_2 ,

$$\begin{aligned} E_\varepsilon W_{n,n,v}(\lambda_{n,v}) &\leq 9\lambda_{n,v}^2 + E_{X,\varepsilon}(W_{n,2,v}(2\lambda_{n,v})), \\ &\leq 9\lambda_{n,v}^2 + Q_{n,v}(2\lambda_{n,v}), \text{ thanks to Lemma 3.5.5, page 88,} \\ &\leq 9\lambda_{n,v}^2 + 4\Delta\lambda_{n,v}^2. \end{aligned}$$

Applying once again Lemma 3.5.8, page 89, we get that

$$W_{n,n,v}(\lambda_{n,v}) \leq E_\varepsilon W_{n,n,v}(\lambda_{n,v}) + \lambda_{n,v}^2 \leq (10 + 4\Delta)\lambda_{n,v}^2.$$

This ends the proof of the lemma by taking $\kappa = 10 + 4\Delta$. □

3.6 Proof of Corollary 3.3.1

According to Theorem 3.3.1 we have with high probability,

$$\|\hat{f} - m\|_n^2 \leq C \inf_{f \in \mathcal{F}} \left\{ \|m - f\|_n^2 + \sum_{v \in \mathcal{S}_f} (\mu_v + \gamma_v^2) \right\}. \quad (3.81)$$

Besides, for all $K > 0$,

$$\|\hat{f} - m\|_2^2 \leq (1 + K)\|\hat{f} - f\|_2^2 + \left(1 + \frac{1}{K}\right)\|m - f\|_2^2. \quad (3.82)$$

We consider once again two cases defined in page 87.

Case 1: $\|\hat{f} - f\|_2 \leq \|\hat{f} - f\|_n$,

In this case Equation (3.82) gives,

$$\|\hat{f} - m\|_2^2 \leq (1 + K)\|\hat{f} - f\|_n^2 + \left(1 + \frac{1}{K}\right)\|m - f\|_2^2.$$

Then, using Equations (3.55) and (3.81) we obtain the result.

Case 2: $\|\hat{f} - f\|_2 \geq \|\hat{f} - f\|_n$,

Apply Lemma 3.5.4 (page 87) and conclude that conditioning on the events \mathcal{T} and \mathcal{A} , defined by (3.47) and (3.52), then $\hat{f} - f$ belongs to $\mathcal{G}(f)$ defined in Lemma 3.5.4. Now, conditioning on the event \mathcal{C} we get the result as in Case 1 since,

$$\|\hat{f} - f\|_2 \leq \sqrt{2}\|\hat{f} - f\|_n.$$

□

Appendix

3.A Proofs of Section 3.4.2

3.A.1 Proof of Remark 3.4.1

From Lemma 3.4.1 we have $U_{\tilde{\alpha}}(u) \subset (u^{1/2}B_2 + u^{1/\tilde{\alpha}}B_{\tilde{\alpha}})$. It suffices to show that $(u^{1/2}B_2 + u^{1/\tilde{\alpha}}B_{\tilde{\alpha}}) \subset 2 \times \max(u^{1/2}, u^{1/\tilde{\alpha}})B_2$.

Consider $x \in u^{1/2}B_2 + u^{1/\tilde{\alpha}}B_{\tilde{\alpha}}$, $x = y + z$ with $y \in u^{1/2}B_2$, means $\sum_{i=1}^n y_i^2 \leq u$, and $z \in u^{1/\tilde{\alpha}}B_{\tilde{\alpha}}$, means $\sum_{i=1}^n z_i^{\tilde{\alpha}} \leq u$. Moreover, we know that $\|x\| \leq \|y\| + \|z\|$ which leads to $\|x\| \leq u^{1/2} + u^{1/\tilde{\alpha}}$ and $\|x\|^2 \leq 2(u + u^{2/\tilde{\alpha}}) \leq 4 \times \max(u, u^{2/\tilde{\alpha}})$. \square

3.A.2 Proof of Corollary 3.4.1

From Equation (3.35) we have $N(2 \times \max(M^{1/2}, M^{1/\tilde{\alpha}}), T, \|\cdot\|) \leq \exp(KM)$.

Using this on sT for $s > 0$ we have $sM = E_Z \sup_{t' \in sT} \sum_{i=1}^n t'_i Z_i$ and,

$$N(2 \times \max((sM)^{1/2}, (sM)^{1/\tilde{\alpha}}), sT, \|\cdot\|) \leq \exp(KsM).$$

Moreover,

$$N(2 \times \max((sM)^{1/2}, (sM)^{1/\tilde{\alpha}}), sT, \|\cdot\|) = N\left(\frac{2}{s} \times \max((sM)^{1/2}, (sM)^{1/\tilde{\alpha}}), T, \|\cdot\|\right),$$

since for all $t_1, t_2 \in T$ and some constant C , $\|st_1 - st_2\| \leq C$ is equivalent to $\|t_1 - t_2\| \leq C/s$.

We obtain then,

$$N\left(\frac{2}{s} \times \max((sM)^{1/2}, (sM)^{1/\tilde{\alpha}}), T, \|\cdot\|\right) \leq \exp(KsM).$$

As in Remark 3.4.2 for $u = sM$ we consider two following cases (recall that $1 < \tilde{\alpha} < 2$):

(i) If $sM \leq 1$ we have $(sM)^{1/\tilde{\alpha}} \leq (sM)^{1/2}$ and so,

$$N\left(2\left(\frac{M}{s}\right)^{1/2}, T, \|\cdot\|\right) \leq \exp(KsM).$$

Take $\delta = 2(M/s)^{1/2}$ and thus $s = 4M/\delta^2$. Moreover, $sM \leq 1$ (i.e. $(4M/\delta^2) \times M \leq 1$) and so $\delta \geq 2M$. Finally, we obtain in this case:

$$\forall \delta \geq 2M, \log N(\delta, T, \|\cdot\|) \leq K\left(\frac{2M}{\delta}\right)^2.$$

(ii) If $sM \geq 1$ we have $(sM)^{1/2} \leq (sM)^{1/\tilde{\alpha}}$ and so,

$$N\left(\frac{2}{s}(sM)^{1/\tilde{\alpha}}, T, \|\cdot\|\right) \leq \exp(KsM).$$

Take $\delta = (2/s)(sM)^{1/\tilde{\alpha}}$ and thus $s = (2/\delta)^{\tilde{\alpha}/(\tilde{\alpha}-1)}M^{1/(\tilde{\alpha}-1)}$. Moreover, $sM \geq 1$ (i.e. $(2M/\delta)^{\tilde{\alpha}/(\tilde{\alpha}-1)} \geq 1$) and so $0 < \delta \leq 2M$. Finally, we obtain in this case:

$$\forall 0 < \delta \leq 2M, \log N(\delta, T, \|\cdot\|) \leq K\left(\frac{2M}{\delta}\right)^{\tilde{\alpha}/(\tilde{\alpha}-1)} = K\left(\frac{2M}{\delta}\right)^\alpha.$$

□

3.B Proofs of Section 3.4.3

3.B.1 Proof of Lemma 3.4.2

In order to prove this Lemma it suffices to show that $\Pi_\alpha \in \mathcal{M}(m, \rho^2)$ for some m . To do so, we use Example 3.4.1.

First show $d \log \Pi_\alpha([t, \infty))/dt \leq -t/\rho^2$:

We know that

$$\frac{d}{dt} \log \Pi_\alpha([t, \infty)) = \frac{d}{dt} \log(1 - \Pi_\alpha((-\infty, t])) = -\frac{\pi_\alpha(t)}{1 - \Pi_\alpha((-\infty, t])} = -\frac{\pi_\alpha(t)}{\Pi_\alpha([t, \infty))}.$$

For all $t > 0$ we have,

$$\Pi_\alpha([t, \infty)) = \int_t^\infty a_\alpha \exp(-|x|^\alpha) dx = \int_t^\infty a_\alpha \exp(-x^\alpha) dx.$$

Take $x = u^{1/\alpha}$, so $dx = (1/\alpha)u^{(1/\alpha)-1} du$, and

$$\begin{aligned} \Pi_\alpha([t, \infty)) &= \int_{t^\alpha}^\infty \frac{a_\alpha}{\alpha} u^{(1/\alpha)-1} \exp(-u) du, \\ &= \frac{a_\alpha}{\alpha} \Gamma\left(\frac{1}{\alpha}, t^\alpha\right), \end{aligned}$$

where $\Gamma(\frac{1}{\alpha}, t^\alpha)$ is incomplete gamma function. Moreover, for $s \in \mathbb{R}$ as $x \rightarrow \infty$,

$$\frac{\Gamma(s, x)}{x^{s-1} \exp(-x)} \rightarrow 1.$$

Therefore,

$$\Pi_\alpha([t, \infty)) = \frac{a_\alpha}{\alpha} t^{1-\alpha} \exp(-t^\alpha).$$

Since $t > 0$ so $\pi_\alpha(t) = a_\alpha \exp(-t^\alpha)$, and

$$\frac{d}{dt} \log \Pi_\alpha([t, \infty)) = -\frac{\alpha a_\alpha \exp(-t^\alpha)}{a_\alpha t^{1-\alpha} \exp(-t^\alpha)} = -\alpha t^{\alpha-1}.$$

The inequality $-\alpha t^{\alpha-1} \leq -t/\rho^2$ (i.e. $t^{\alpha-2} \geq 1/\alpha\rho^2$) holds for all $\alpha > 2$ and $t \geq (1/\alpha\rho^2)^{1/(\alpha-2)}$.

Second show $d \log \Pi_\alpha((-\infty, -t])/dt \leq -t/\rho^2$:

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
110 in the regression model with non-Gaussian and non-bounded error**

The probability distribution Π_α is symmetric, therefore $\Pi_\alpha((-\infty, -t]) = \Pi_\alpha([t, \infty))$, and

$$\frac{d}{dt} \log \Pi_\alpha((-\infty, -t]) = \frac{d}{dt} \log \Pi_\alpha([t, \infty)) = -\alpha t^{\alpha-1},$$

which is smaller than $-t/\rho^2$ if $\alpha > 2$ and $t \geq (1/\alpha\rho^2)^{1/(\alpha-2)}$.

Take $m = (1/\alpha\rho^2)^{1/(\alpha-2)}$, then for $x \geq m$, Π_α verifies the Equations (3.39). That is $\Pi_\alpha \in \mathcal{M}((1/\alpha\rho^2)^{1/(\alpha-2)}, \rho^2)$. □

3.B.2 Proof of Remark 3.4.3

If $\alpha = 2$, according to the Laplace transform of the Gaussian function we have

$$E(\exp(s|Z|)) = 2a_\alpha\sqrt{\pi} \exp\left(\frac{s^2}{4}\right). \quad (3.83)$$

If $\alpha > 2$ we have,

$$E(\exp(s|Z|)) = \int_{-\infty}^{+\infty} \exp(s|z|)a_\alpha \exp(-|z|^\alpha)dz = 2a_\alpha\mathcal{S},$$

where

$$\mathcal{S} = \int_0^{+\infty} \exp(sz - z^\alpha)dz = \underbrace{\int_0^1 \exp(sz - z^\alpha)dz}_{\mathcal{S}_1} + \underbrace{\int_1^{+\infty} \exp(sz - z^\alpha)dz}_{\mathcal{S}_2}.$$

For $z \in [0, 1]$ we have $\exp(-z^\alpha) \leq 1$ and so

$$\mathcal{S}_1 \leq \int_0^1 \exp(sz)dz = \frac{\exp(s) - 1}{s}.$$

For $z \geq 1$ we have $\exp(z^2 - z^\alpha) < 1$ and so

$$\mathcal{S}_2 = \int_1^{+\infty} \exp(sz - z^2 + z^2 - z^\alpha)dz < \int_1^{+\infty} \exp(sz - z^2)dz < \sqrt{\pi} \exp\left(\frac{s^2}{4}\right),$$

where the last inequality is obtained using Equation (3.83). Finally, we obtain

$$\mathcal{S} < \frac{\exp(s) - 1}{s} + \sqrt{\pi} \exp\left(\frac{s^2}{4}\right),$$

and therefore

$$E(\exp(s|Z|)) < 2a_\alpha\left(\frac{\exp(s) - 1}{s} + \sqrt{\pi} \exp\left(\frac{s^2}{4}\right)\right).$$

□

3.B.3 Proof of Corollary 3.4.2

We suppose that the inequality (3.41) holds and we want to find an upper bound for $P(|\phi(Z) - E(\phi(Z))| \geq u)$. Using the Markov's inequality we have,

$$\begin{aligned} P(|\phi(Z) - E(\phi(Z))| > u) &= P(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2) > \exp(\lambda u^2)), \\ &\leq \exp(-\lambda u^2) E(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)). \end{aligned} \quad (3.84)$$

To demonstrate the result of the Theorem, it suffices to find an upper bound for the following quantity

$$E(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)).$$

Let Z_1 and Z_2 be two independent random variables distributed with the same law, then for all $u' > 0$ we have:

$$P(|Z_1 - Z_2| > u') \leq P(|Z_1 - M(\phi(Z_1))| > \frac{u'}{2}) + P(|Z_2 - M(\phi(Z_2))| > \frac{u'}{2}). \quad (3.85)$$

Furthermore, for all convex function ψ we have:

$$E(\psi(Z_1 - E(Z_1))) = \int \psi\left(\int (z_1 - z_2) dP(z_2)\right) dP(z_1).$$

Applying the Jensen's inequality we obtain then,

$$\begin{aligned} E(\psi(Z_1 - E(Z_1))) &\leq \int \left(\int \psi(z_1 - z_2) dP(z_2)\right) dP(z_1), \\ &\leq E(\psi(Z_1 - Z_2)). \end{aligned} \quad (3.86)$$

Set $\psi(t) = \exp(\lambda t^2)$ for $\lambda > 0$, $Z_1 = \phi(Z)$ and $Z_2 = \phi(Z')$. Since $\psi(t)$ is convex, then Equation (3.86) gives:

$$E(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)) \leq E(\exp(\lambda(\phi(Z) - \phi(Z'))^2)). \quad (3.87)$$

For all non-negative random variables Z we have $E(Z) = \int_{[0, \infty)} P(Z \geq z) dz$. So, we obtain from Equation (3.87):

$$E(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)) \leq \int_0^\infty P(\exp(\lambda(\phi(Z) - \phi(Z'))^2) > t) dt.$$

Using Equation (3.85) and simple calculations leads to:

$$\begin{aligned} E(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)) &\leq \int_1^\infty P(|\phi(Z) - \phi(Z')| > \sqrt{\frac{\log(t)}{\lambda}}) dt, \\ &\leq 2 \int_1^\infty P(|\phi(Z) - M(\phi(Z))| > \frac{1}{2} \sqrt{\frac{\log(t)}{\lambda}}) dt. \end{aligned} \quad (3.88)$$

**Chapter 3. Risk upper bounds for RKHS ridge group sparse estimator
112 in the regression model with non-Gaussian and non-bounded error**

In this step we can use the result in Equation (3.41), from which we obtain:

$$E\left(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)\right) \leq 2B \int_1^\infty \exp\left(-\frac{\log(t)}{4\lambda A}\right) dt, \quad (3.89)$$

and, therefore,

$$E\left(\exp(\lambda|\phi(Z) - E(\phi(Z))|^2)\right) \leq \frac{8\lambda AB}{1 - 4\lambda A}, \quad \forall \lambda < \frac{1}{4A}. \quad (3.90)$$

The proof is complete by taking $\lambda = 1/8A$.

□

Estimate the Hoeffding decomposition of a complex model by solving RKHS ridge group sparse optimization problem

Abstract

We propose an R package, called **RKHSMetaMod**, that implements a procedure for estimating a meta-model of a complex model m . The meta-model approximates the Hoeffding decomposition of m and allows to perform sensitivity analysis on it. It belongs to a reproducing kernel Hilbert space that is constructed as a direct sum of Hilbert spaces. The estimator of the meta-model is the solution of a penalized least-squares minimization with the sum of the Hilbert norm and the empirical L^2 -norm. This procedure, called RKHS ridge group sparse, allows both to select and estimate the terms in the Hoeffding decomposition, and therefore, to select and estimate the Sobol indices that are non-zero. This package provides an interface from R statistical computing environment to the C++ libraries **Eigen** and **GSL**. In order to speed up the execution time and optimize the storage memory, except for a function that is written in R, all of the functions of **RKHSMetaMod** package are written using the efficient C++ libraries through **RcppEigen** and **RcppGSL** packages. These functions are then interfaced in the R environment in order to propose an user friendly package.

Keywords: meta-model, Hoeffding decomposition, ridge group sparse, reproducing kernel Hilbert space, Sobol indices.

4.1 Introduction

Let us consider a phenomenon described by a model m depending on d input variables $X = (X_1, \dots, X_d)$. This model m from \mathbb{R}^d to \mathbb{R} may be a known model that can be calculated in all points of X , or it may be an unknown regression model defined as follows:

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0, \quad (4.1)$$

where the error ε is assumed to be centered with a finite variance, i.e. $E(\varepsilon) = 0$ and $\text{var}(\varepsilon) < \infty$.

The components of X are independent and have a known law $P_X = \prod_{a=1}^d P_{X_a}$ on \mathcal{X} , a subset of \mathbb{R}^d . The number d of components of X may be large. The model m may present high complexity as strong non-linearities and high order interaction effects, and it is assumed to be square-integrable, i.e. $m \in L^2(\mathcal{X}, P_X)$.

Based on n data points $\{(X_i, Y_i)\}_{i=1}^n$, a meta-model that approximates the Hoeffding decomposition of m is estimated. This meta-model belongs to a reproducing kernel Hilbert space (RKHS), which is constructed as a direct sum of Hilbert spaces leading to an additive decomposition including the variables and interactions between them (Durrande et al. (2013)). The estimator of the meta-model is calculated by minimizing a least-squares criterion penalized by the sum of two penalty terms: the Hilbert norm and the empirical norm (Huet and Taupin (2017)). This procedure allows to select the subsets of variables X that contribute to predict Y . The estimated meta-model is used to perform sensitivity analysis, and so allows to determine the influence of each variable and groups of them on the output variable Y .

In the classical framework of sensitivity analysis $m(X)$ is calculable in all points of X . In this framework, one may use the method of Sobol (1993) for variance-based methods of global sensitivity analysis in order to perform sensitivity analysis on m . Let us briefly recall this method.

Let \mathcal{P} be the set of all subsets of $\{1, \dots, d\}$ with dimension 1 to d . For all $X \in \mathcal{X}$ and $v \in \mathcal{P}$, let X_v be the vector with components X_a for all $a \in v$. For a set A let $|A|$ be its cardinality, and for all $v \in \mathcal{P}$, let $m_v : \mathbb{R}^{|v|} \rightarrow \mathbb{R}$ be a function of X_v .

The independency between the components of X allows to write the function m according to its Hoeffding decomposition (Sobol (1993), van der Vaart (1998)):

$$m(X) = m_0 + \sum_{v \in \mathcal{P}} m_v(X_v), \quad (4.2)$$

where m_0 is known as constant term, when $|v| = 1$ the functions m_v are known as main effects, when $|v| = 2$, i.e. $v = \{a, a'\}$ and $a \neq a'$, the functions m_v are known as second order interactions, and so on.

This decomposition (4.2) is unique, all the terms m_v , $v \in \mathcal{P}$ are centered, and they are orthogonal with respect to $L^2(\mathcal{X}, P_X)$. The function m as well as all the functions m_v in Equation (4.2) are square-integrable. As any two terms of decomposition (4.2) are orthogonal, by squaring (4.2) and integrating it with respect to the distribution of X , a decomposition of the variance of $m(X)$ is obtained as follows:

$$\text{var}(m(X)) = \sum_{v \in \mathcal{P}} \text{var}(m_v(X_v)). \quad (4.3)$$

For any group of variables X_v , $v \in \mathcal{P}$, the Sobol indices are defined by:

$$S_v = \frac{\text{var}(m_v(X_v))}{\text{var}(m(X))}. \quad (4.4)$$

For each v , S_v expresses the fraction of variance of $m(X)$ explained by X_v .

For all $v \in \mathcal{P}$, when $|v| = 1$, the S_v 's are referred to as the first order indices. When $|v| = 2$, i.e. $v = \{a, a'\}$ and $a \neq a'$, they are referred to as the second order indices or the interaction indices of order two (between X_a and $X_{a'}$). And the same holds for $|v| > 2$.

The total number of the Sobol indices to be calculated is equal to $|\mathcal{P}| = 2^d - 1$, which raises exponentially with the number d of the input variables. When d is large, the evaluation of all the indices can be too computationally demanding and even not reachable. For this reason, only the indices of order not higher than two are calculated in practice. However, only first and second order indices may not provide a good information on the model sensitivities. In order to provide a better information on the model sensitivities, Homma and Saltelli (1996) proposed to calculate the first order and the total indices defined as follows:

Let $\mathcal{P}_a \subset \mathcal{P}$ be the set of all the subsets of $\{1, \dots, d\}$ including a , then

$$S_{T_a} = \sum_{v \in \mathcal{P}_a} S_v.$$

For all $a \in \{1, \dots, d\}$, S_{T_a} denotes the total effect of X_a . It expresses the fraction of variance of $m(X)$ explained by X_a alone and all the interactions of it with the other variables.

The total indices allow to rank the input variables with respect to the amount of their effect on the output variable. However, they do not provide complete information on the model sensitivities as do all the Sobol indices.

The classical computation of the Sobol indices is based on the Monte Carlo methods (see for example: Sobol (1993) for the main effect and interaction indices, and Saltelli (2002) for the main effect and total indices). For models that are expensive to evaluate, the Monte Carlo methods lead to high computational burden. Moreover, in the case where d is large, m is complex and the calculation of the variances (see Equation (4.3)) is numerically complicated or not possible, as in the case where the model m is unknown, the methods described above are not applicable.

Another method is to approximate m by a simplified model, called a meta-model, which is much faster to evaluate and to perform sensitivity analysis on it. A meta-model provides additional information than just scalar indices. It provides the approximations of the Sobol indices of m at a lower computational cost, and also a deeper view of the input variable's effects on the model output.

Among the meta-modelling methods proposed in the literature, the expansion based on polynomial Chaos (Wiener (1938), Schoutens (2000)) can be used to approximate the Hoeffding decomposition of m (Sudret (2008)).

The principle of the polynomial Chaos is to project m onto a basis of orthonormal polynomials. The polynomial Chaos expansion of m is written as (Soize and Ghanem (2004)):

$$m(X) = \sum_{j=0}^{\infty} h_j \phi_j(X), \quad (4.5)$$

where $\{h_j\}_{j=0}^\infty$ are the coefficients, and $\{\phi_j\}_{j=0}^\infty$ are multivariate orthonormal polynomials associated with X that are determined according to the distribution of the components of X . In practice, expansion (4.5) shall be truncated for computational purposes, and the model m may be approximated by:

$$m(X) \approx \sum_{j=0}^{v_{max}} h_j \phi_j(X),$$

where v_{max} is determined using a *truncation scheme*. In this approach, the Sobol indices are obtained by summing up the squares of the suitable coefficients.

Blatman and Sudret (2011) proposed a method for truncating the polynomial Chaos expansion and an algorithm based on least angle regression for selecting the terms in the expansion.

In this method, according to the distribution of the components of X a unique family of orthonormal polynomials $\{\phi_j\}_{j=0}^\infty$ is determined. However, this family may not be necessarily the best functional basis to approximate m well.

Another method to construct meta-models is the Gaussian Process (GP) modelling (Welch et al. (1992), Oakley and O'Hagan (2004), Kleijnen (2007, 2009), Marrel et al. (2009), Durrande et al. (2012), Le Gratiet et al. (2014)). The principle is to consider that the prior knowledge about the function $m(X)$, can be modelled by a GP $\mathcal{Z}(X)$ with a mean $m_{\mathcal{Z}}(X)$ and a covariance kernel $k_{\mathcal{Z}}(X, X')$. To perform sensitivity analysis from a GP model one may replace the true model $m(X)$ with the mean of the conditional GP, and deduce the Sobol indices from it.

A review on the meta-modelling based on polynomial Chaos and GP is presented in Le Gratiet et al. (2017).

Durrande et al. (2013) considered a class of functional approximation methods similar to the GP and obtained a meta-model that satisfies the properties of the Hoeffding decomposition. They proposed to approximate m by functions belonging to a RKHS \mathcal{H} which is constructed as a direct sum of Hilbert spaces such that the projection of m onto \mathcal{H} , denoted f^* , is an approximation of the Hoeffding decomposition of m .

The function f^* is defined as the minimizer over the functions $f \in \mathcal{H}$ of the following criterion,

$$E_X(m(X) - f(X))^2.$$

Let $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ be the inner product in \mathcal{H} , let also k and k_v be the reproducing kernels associated with the RKHS \mathcal{H} and the RKHS \mathcal{H}_v , respectively. The properties of the RKHS \mathcal{H} insures that any function $f \in \mathcal{H}$, $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ is written as the following decomposition:

$$f(X) = \langle f, k(X, \cdot) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \quad (4.6)$$

where f_0 is a constant, and $f_v : \mathbb{R}^{|v|} \rightarrow \mathbb{R}$ is defined by,

$$f_v(X) = \langle f, k_v(X, \cdot) \rangle_{\mathcal{H}}.$$

For all $v \in \mathcal{P}$, the functions $f_v(X_v)$ are centered and for all $v \neq v'$, the functions $f_v(X_v)$ and $f_{v'}(X_{v'})$ are orthogonal with respect to $L^2(\mathcal{X}, P_X)$. So the decomposition of the function f presented in Equation (4.6) is its Hoeffding decomposition. As the function f^* belongs to the RKHS \mathcal{H} , it is decomposed as its Hoeffding decomposition:

$$f^* = f_0^* + \sum_{v \in \mathcal{P}} f_v^*, \quad (4.7)$$

and each function f_v^* approximates the function m_v in Equation (4.2). In the decomposition (4.7), we have $|\mathcal{P}|$ terms f_v^* to be estimated. The cardinality of \mathcal{P} is equal to $2^d - 1$ which may be huge since it raises very quickly by increasing d . In order to deal with this problem, in the regression framework, one may estimate f^* by a sparse meta-model $\hat{f} \in \mathcal{H}$. To this purpose, the estimation of f^* is done on the basis of n observations by minimizing a least-squares criterion suitably penalized in order to deal both with the non-parametric nature of the problem, and with the possibly large number of functions that have to be estimated.

Note that, in the classical framework of sensitivity analysis, where $m(X)$ is calculable in all points X , one may calculate a sparse approximation of f^* using least-squares penalized criterion as it is done in the non-parametric regression framework.

In order to obtain a sparse solution of a minimization problem, the penalty function should enforce the sparsity. There exists various ways of enforcing sparsity for a minimization (maximization) problem, see for example [Hastie et al. \(2015\)](#) for a review. Some methods, such as the Sparse Additive Models (SpAM) procedure ([Ravikumar et al. \(2009\)](#), [Liu et al. \(2009\)](#)) are based on a combination of the l_1 -norm with the empirical L^2 -norm,

$$\|f\|_{n,1} = \sum_{a=1}^d \|f_a\|_n,$$

where

$$\|f_a\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_a^2(X_{ai}),$$

is the squared empirical L^2 -norm of the univariate function f_a . The Component Selection and Smoothing Operator (COSSO) method developed by [Lin and Zhang \(2006\)](#) enforces sparsity using a combination of the l_1 -norm with the Hilbert norm,

$$\|f\|_{\mathcal{H},1} = \sum_{a=1}^d \|f_a\|_{\mathcal{H}_a}.$$

Instead of focusing on only one penalty term, one may consider a more general family of estimators, called *doubly penalized estimator*, that is obtained by minimizing a criterion penalized by the sum of two penalty terms. [Raskutti et al. \(2009, 2012\)](#) proposed a *doubly penalized estimator* which is the solution of the minimization of a least-squares criterion penalized by the sum of a sparsity penalty term and a

combination of the l_1 -norm with the Hilbert norm,

$$\gamma \|f\|_{n,1} + \mu \|f\|_{\mathcal{H},1}, \quad (4.8)$$

where $\gamma, \mu \in \mathbb{R}$ are the tuning parameters that should be suitably chosen.

Meier et al. (2009) proposed a related family of estimators, based on the penalization with the empirical L^2 -norm. Their penalty function is the sum of the sparsity penalty term, $\|f\|_{n,1}$, and a smoothness penalty term.

Huet and Taupin (2017) considered the same approximation functional spaces as Durrande et al. (2013), and obtained a *doubly penalized estimator* of a meta-model which approximates the Hoeffding decomposition of m . Their estimator is the solution of least-squares minimization penalized by the penalty function defined in Equation (4.8) adapted to the multivariate setting,

$$\gamma \|f\|_n + \mu \|f\|_{\mathcal{H}}, \quad (4.9)$$

with

$$\|f\|_n = \sum_{v \in \mathcal{P}} \|f_v\|_n, \text{ and } \|f\|_{\mathcal{H}} = \sum_{v \in \mathcal{P}} \|f_v\|_{\mathcal{H}_v}.$$

This procedure, called RKHS ridge group sparse, estimates the groups v that are suitable for predicting f^* , and the relationship between f_v^* and X_v for each group. The obtained estimator, called RKHS meta-model, is used then to estimate the Sobol indices of m . This approach makes it possible to estimate the Sobol indices for all groups in the support of the RKHS meta-model, including the interactions of possibly high order, a point known to be difficult in practice.

In this Chapter, an R package, called **RKHSMetaMod**, that implements the RKHS ridge group sparse procedure is proposed. This package deals with the input variables $X = (X_1, \dots, X_d)$ that are independent and uniformly distributed on $\mathcal{X} = [0, 1]^d$, i.e. $X \sim P_X = P_1 \times \dots \times P_d$, with P_a , $a = 1, \dots, d$ being the uniform law on the interval $[0, 1]$. It allows to:

- (1) calculate reproducing kernels and their associated Gram matrices (see Section 4.3.1),
- (2) implement the RKHS ridge group sparse procedure and a special case of it called the RKHS group lasso procedure, i.e. when $\gamma = 0$ in the penalty function (4.9), in order to estimate the terms f_v^* in the Hoeffding decomposition of f^* leading to an estimation of the function m (see Section 4.3.2),
- (3) choose the tuning parameters μ and γ (see Equation (4.9)), using a procedure that leads to obtain the *best* RKHS meta-model in terms of the prediction quality,
- (4) estimate the Sobol indices of the function m (see Section 4.2.4).

To the best of our knowledge, there is no other package available that implements the RKHS ridge group sparse procedure. The **RKHSMetaMod** package is dedicated to the meta-model estimation on the RKHS \mathcal{H} . The convex optimization algorithms used in this package are adapted to take into account the problem of high dimensionality in this context. This package is available from the Comprehensive R Archive Network (CRAN) at <https://cran.r-project.org/web/packages/RKHSMetaMod/>.

The organization of this Chapter is as follows: In Section 4.2, the estimation method is described. In Section 4.3, the algorithms used in the **RKHSMetaMod** package to obtain the RKHS meta-model are detailed. In Section 4.4, an overview of the package functions as well as a brief documentation of them are given. In Section 4.5, a simulation study to validate the performances of the **RKHSMetaMod** package functions is given.

4.2 Estimation method

In Section 4.2.1, the RKHS ridge group sparse and the RKHS group lasso procedures are presented. In Section 4.2.2, the method of [Durrande et al. \(2013\)](#) to construct the RKHS \mathcal{H} is recalled. The strategy of choosing the tuning parameters in the RKHS ridge group sparse algorithm is detailed in Section 4.2.3, and in Section 4.2.4, the calculation of the empirical Sobol indices of the RKHS meta-model is described.

4.2.1 RKHS ridge group sparse and RKHS group lasso procedures

Let denote by n , the number of observations. The dataset consists of a vector of n observations $Y = (Y_1, \dots, Y_n)$, and a $n \times d$ matrix of features X with components,

$$(X_{ai}, i = 1, \dots, n, a = 1, \dots, d) \in \mathbb{R}^{n \times d}.$$

For some tuning parameters $\gamma_v, \mu_v, v \in \mathcal{P}$, the RKHS ridge group sparse criterion is defined by,

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \left(Y_i - f_0 - \sum_{v \in \mathcal{P}} f_v(X_{vi}) \right)^2 + \sum_{v \in \mathcal{P}} \gamma_v \|f_v\|_n + \sum_{v \in \mathcal{P}} \mu_v \|f_v\|_{\mathcal{H}_v}, \quad (4.10)$$

where X_v represents the matrix of variables corresponding to the v -th group,

$$X_v = (X_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|},$$

and where $\|f_v\|_n$ is the empirical L^2 -norm of f_v defined by the sample $\{X_{vi}\}_{i=1}^n$ as,

$$\|f_v\|_n^2 = \frac{1}{n} \sum_{i=1}^n f_v^2(X_{vi}).$$

The penalty function in the criterion (4.10) is the sum of the Hilbert norm and the empirical norm, which allows to select few terms in the additive decomposition

of f over sets $v \in \mathcal{P}$. Moreover, the Hilbert norm favours the smoothness of the estimated f_v , $v \in \mathcal{P}$.

Let us define the set of functions,

$$\mathcal{F} = \left\{ f : f = f_0 + \sum_{v \in \mathcal{P}} f_v, \text{ with } f_v \in \mathcal{H}_v, \text{ and } \|f_v\|_{\mathcal{H}_v} \leq r_v, r_v > 0 \right\}.$$

Then the RKHS meta-model is defined by,

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f). \quad (4.11)$$

According to the Representer Theorem (Kimeldorf and Wahba (1970)) the non-parametric functional minimization problem described above is equivalent to a parametric minimization problem. Indeed, the solution of the minimization problem (4.11) belonging to the RKHS \mathcal{H} is written as $f = f_0 + \sum_{v \in \mathcal{P}} f_v$, where for some matrix $\theta = (\theta_{vi}, i = 1, \dots, n, v \in \mathcal{P}) \in \mathbb{R}^{n \times |\mathcal{P}|}$ we have for all $v \in \mathcal{P}$,

$$f_v(\cdot) = \sum_{i=1}^n \theta_{vi} k_v(X_{vi}, \cdot).$$

Let $\|\cdot\|$ be the Euclidean norm in \mathbb{R}^n , and for each $v \in \mathcal{P}$, let K_v be the $n \times n$ Gram matrix associated with the kernel $k_v(\cdot, \cdot)$, i.e. $(K_v)_{i,i'} = k_v(X_{vi}, X_{vi'})$. Let also $K_v^{1/2}$ be the matrix that satisfies $t(K_v^{1/2})K_v^{1/2} = K_v$, and let \hat{f}_0 and $\hat{\theta}$ be the minimizers of the following penalized least-squares criterion:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \sum_{v \in \mathcal{P}} \gamma_v \|K_v \theta_v\| + n \sum_{v \in \mathcal{P}} \mu_v \|K_v^{1/2} \theta_v\|.$$

Then the estimator \hat{f} defined in Equation (4.11) satisfies,

$$\hat{f}(X) = \hat{f}_0 + \sum_{v \in \mathcal{P}} \hat{f}_v(X_v) \text{ with } \hat{f}_v(X_v) = \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v).$$

Remark 4.2.1 *The constraint $\|f_v\|_{\mathcal{H}_v} \leq r_v$ is not taken into account in the parametric minimization problem. This constraint is crucial for theoretical properties but the value of r_v is unknown and has no practical usefulness.*

For each $v \in \mathcal{P}$, let γ'_v and μ'_v be the weights that are chosen suitably. We define,

$$\gamma_v = \gamma \times \gamma'_v \text{ and } \mu_v = \mu \times \mu'_v \text{ with } \gamma, \mu \in \mathbb{R}^+.$$

Remark 4.2.2 *This formulation simplify the choice of the tuning parameters, since instead of tuning the parameters γ_v and μ_v for all $v \in \mathcal{P}$, only two parameters γ and μ are tuned. Moreover, the weights γ'_v and μ'_v , $v \in \mathcal{P}$, may be of interest in applications. For example, one can take weights that increase with the cardinal of v in order to favour effects with small interaction order between variables.*

For the sake of simplicity, in the rest of this Chapter for all $v \in \mathcal{P}$ the weights γ'_v and μ'_v are assumed to be setted as 1, and the RKHS ridge group sparse criterion is then expressed as follows:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \gamma \sum_{v \in \mathcal{P}} \|K_v \theta_v\| + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\|. \quad (4.12)$$

By considering only the second part of the penalty function in the RKHS ridge group sparse criterion (4.12), i.e. by setting $\gamma = 0$, the RKHS group lasso criterion is obtained as follows:

$$C_g(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\|, \quad (4.13)$$

which is a group lasso criterion (Yuan and Lin (2006)) up to a scale transformation.

In the **RKHSMetaMod** package, the RKHS ridge group sparse algorithm is initialized using the solutions obtained by solving the RKHS group lasso algorithm. Indeed, the penalty function in the RKHS group lasso criterion (4.13) insures the sparsity in the solution. Therefore, for a given value of μ , by implementing the RKHS group lasso algorithm (see Section 4.3.2.1), a RKHS meta-model with few terms in its additive decomposition is obtained. The support and the coefficients of a RKHS meta-model which is obtained by implementing RKHS group lasso algorithm will be denoted by $\widehat{S}_{f_{\text{Group Lasso}}}$ and $\widehat{\theta}_{\text{Group Lasso}}$, respectively.

From now on the tuning parameter in the RKHS group lasso criterion will be denoted by:

$$\mu_g = \sqrt{n} \mu. \quad (4.14)$$

4.2.2 RKHS construction

We consider a RKHS \mathcal{H} that is constructed as a direct sum of Hilbert spaces, and that is associated with a so-called ANOVA kernel. The ANOVA kernel is defined in order to obtain the analytical expression of the terms of the Hoeffding decomposition of the functions of \mathcal{H} . Therefore, any function f in \mathcal{H} is a candidate to approximate the Hoeffding decomposition of m . The construction of the RKHS \mathcal{H} has been proposed by Durrande et al. (2013) that we recall briefly in the following.

Let $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ be a subset of \mathbb{R}^d . For each $a \in \{1, \dots, d\}$, we choose a RKHS \mathcal{H}_a and its associated kernel k_a defined on the set $\mathcal{X}_a \subset \mathbb{R}$ such that the two following properties are satisfied:

- (i) $k_a : \mathcal{X}_a \times \mathcal{X}_a \rightarrow \mathbb{R}$ is $P_a \otimes P_a$ measurable,
- (ii) $E_{X_a} \sqrt{k_a(X_a, X_a)} < \infty$.

The property (ii) depends on the kernel k_a , $a = 1, \dots, d$ and the distribution of X_a , $a = 1, \dots, d$. It is not very restrictive since it is satisfied, for example, for any bounded kernel.

The RKHS \mathcal{H}_a can be decomposed as a sum of two orthogonal sub-RKHS,

$$\mathcal{H}_a = \mathcal{H}_{0a} \overset{\perp}{\oplus} \mathcal{H}_{1a},$$

where \mathcal{H}_{0a} is the RKHS of zero mean functions,

$$\mathcal{H}_{0a} = \left\{ f_a \in \mathcal{H}_a : E_{X_a}(f_a(X_a)) = 0 \right\},$$

and \mathcal{H}_{1a} is the RKHS of constant functions,

$$\mathcal{H}_{1a} = \left\{ f_a \in \mathcal{H}_a : f_a(X_a) = C \right\}.$$

The kernel k_{0a} associated with the RKHS \mathcal{H}_{0a} is defined by:

$$k_{0a}(X_a, X'_a) = k_a(X_a, X'_a) - \frac{E_{U \sim P_a}(k_a(X_a, U))E_{U \sim P_a}(k_a(X'_a, U))}{E_{(U, V) \sim P_a \otimes P_a} k_a(U, V)}. \quad (4.15)$$

Let $k_v(X_v, X'_v) = \prod_{a \in v} k_{0a}(X_a, X'_a)$, then the ANOVA kernel $k(., .)$ is defined as follows:

$$k(X, X') = \prod_{a=1}^d \left(1 + k_{0a}(X_a, X'_a) \right) = 1 + \sum_{v \in \mathcal{P}} k_v(X_v, X'_v).$$

For \mathcal{H}_v being the RKHS associated with the kernel k_v , the RKHS associated with the ANOVA kernel is then defined by,

$$\mathcal{H} = \prod_{a=1}^d \left(\mathbb{1} \overset{\perp}{\oplus} \mathcal{H}_{0a} \right) = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v.$$

where \perp denotes the L^2 inner product.

According to this construction, any function $f \in \mathcal{H}$ satisfies decomposition (4.6),

$$f(X) = \langle f, k(X, .) \rangle_{\mathcal{H}} = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v),$$

which is the Hoeffding decomposition of f .

The regularity properties of the RKHS \mathcal{H} constructed as described above, depend on the set of the kernels $(k_a, a = 1, \dots, d)$. This method allows to choose different approximation spaces independently of the distribution of the input variables X_1, \dots, X_d , by choosing different sets of kernels. While as mentioned earlier, in the meta-modelling approach based on polynomial Chaos expansion, according to the distribution of the input variables X_1, \dots, X_d a unique family of orthonormal polynomials $\{\phi_j\}_{j=0}^{\infty}$ is determined. Here, the distribution of the components of X occurs only for the orthogonalization of the spaces $\mathcal{H}_v, v \in \mathcal{P}$, and not in the choice of the RKHS, under the condition that properties (i) and (ii) are satisfied. This is one of the main advantages of this method compared to the method based on the truncated polynomial Chaos expansion where the smoothness of the approximation is handled only by the choice of the truncation (Blatman and Sudret (2011)).

4.2.3 Choice of the tuning parameters

While dealing with an optimization problem, one of the essential steps is to choose appropriately the tuning parameters. To do so,

- ✓ first, a grid of values of the tuning parameters μ and γ is chosen.

Let μ_{\max} be the smallest value of μ_g (see Equation (4.14)), such that the solution to the minimization of the RKHS group lasso problem for all $v \in \mathcal{P}$ is $\theta_v = 0$. We have,

$$\mu_{\max} = \max_v \left(\frac{2}{\sqrt{n}} \|K_v^{1/2}(Y - \bar{Y})\| \right). \quad (4.16)$$

In order to set up the grid of values of μ , one may find μ_{\max} , and then a grid of values of μ is defined as follows:

$$\mu_l = \frac{\mu_{\max}}{(\sqrt{n} \times 2^l)}, \quad l \in \{1, \dots, l_{\max}\}.$$

The grid of values of γ is chosen by the user.

- ✓ next, for the grid of values of μ and γ a sequence of estimators is calculated. Each estimator associated with the pair (μ, γ) in the grid of values of μ and γ , denoted by $\hat{f}_{(\mu, \gamma)}$, is the solution of the RKHS ridge group sparse optimization problem or the RKHS group lasso optimization problem if $\gamma = 0$.
- ✓ finally, the obtained estimators $\hat{f}_{(\mu, \gamma)}$ are evaluated using a testing dataset,

$$\{(Y_i^{\text{test}}, X_i^{\text{test}})\}_{i=1}^{n^{\text{test}}}.$$

The prediction error associated with the estimator $\hat{f}_{(\mu, \gamma)}$ is calculated by,

$$\text{ErrPred}(\mu, \gamma) = \frac{1}{n^{\text{test}}} \sum_{i=1}^{n^{\text{test}}} (Y_i^{\text{test}} - \hat{f}_{(\mu, \gamma)}(X_i^{\text{test}}))^2,$$

where for $S_{\hat{f}}$ being the support of the estimator $\hat{f}_{(\mu, \gamma)}$,

$$\hat{f}_{(\mu, \gamma)}(X^{\text{test}}) = \hat{f}_0 + \sum_{v \in S_{\hat{f}}} \sum_{i=1}^n \hat{\theta}_{vi} k_v(X_{vi}, X_v^{\text{test}}).$$

The pair $(\hat{\mu}, \hat{\gamma})$ with the smallest value of the prediction error is chosen, and the estimator $\hat{f}_{(\hat{\mu}, \hat{\gamma})}$ is considered as the *best* estimator of the function m , in terms of the prediction error.

In the **RKHSMetaMod** package, the algorithm to calculate a sequence of the RKHS meta-models, the value of μ_{\max} , and the prediction error are implemented as **RKHSMetMod**, **mu_max**, and **PredErr** functions, respectively. These functions are described in Section 4.4, and illustrated in Example 4.5.1, Example 4.5.3, and Examples 4.5.1, 4.5.3, 4.5.4, respectively.

4.2.4 Estimation of the Sobol indices

The variance of the function m is estimated by the variance of the estimator \hat{f} . As the estimator \hat{f} belongs to the RKHS \mathcal{H} , it admits the Hoeffding decomposition and,

$$\text{var}(\hat{f}(X)) = \sum_{v \in \mathcal{P}} \text{var}(\hat{f}_v(X_v)),$$

where for all $v \in \mathcal{P}$,

$$\text{var}(\hat{f}_v(X_v)) = E_X(\hat{f}_v^2(X_v)) = \|\hat{f}_v\|_2^2.$$

In order to reduce the computational cost in practice, one may estimate the variances of $\hat{f}_v(X_v)$, $v \in \mathcal{P}$ by their empirical variances.

Let \hat{f}_v be the empirical mean of $\hat{f}_v(X_{vi})$, $i = 1, \dots, n$, then

$$\widehat{\text{var}}(\hat{f}_v(X_v)) = \frac{1}{n-1} \sum_{i=1}^n (\hat{f}_v(X_{vi}) - \hat{f}_v)^2.$$

For the groups v that belong to the support of \hat{f} , the estimators of the Sobol indices of m are defined by,

$$\hat{S}_v = \frac{\widehat{\text{var}}(\hat{f}_v(X_v))}{\sum_{v \in \mathcal{P}} \widehat{\text{var}}(\hat{f}_v(X_v))},$$

and for the groups v that do not belong to the support of \hat{f} , we have $\hat{S}_v = 0$.

In the **RKHSMetaMod** package, the algorithm to calculate the empirical Sobol indices \hat{S}_v , $v \in \mathcal{P}$ is implemented as `SI_emp` function. This function is described in Section 4.4.2 and illustrated in Examples 4.5.1, 4.5.3, 4.5.4.

4.3 Algorithms

The **RKHSMetaMod** package implements two optimization algorithms: the RKHS ridge group sparse (see Algorithm 2), and the RKHS group lasso (see Algorithm 1). These algorithms rely on the Gram matrices K_v , $v \in \mathcal{P}$, that have to be positive definite. Therefore, the first and essential step in this package, is to calculate these matrices and insure their positive definiteness. This step is detailed in an algorithm that is described in Section 4.3.1.

The second step is to estimate the RKHS meta-model. In the **RKHSMetaMod** package, two different objectives based on different procedures are considered in order to calculate this estimator:

1. The RKHS meta-model with the *best* prediction quality:

A sequence of values of the tuning parameters (μ, γ) is considered, and the RKHS meta-models associated with each pair of values of (μ, γ) are calculated. For $\gamma = 0$, the RKHS meta-model is obtained by solving the RKHS

group lasso optimization problem, while for $\gamma \neq 0$ the RKHS ridge group sparse optimization problem is solved to calculate the RKHS meta-model. The obtained meta-models are evaluated by considering a new dataset. The RKHS meta-model with minimum value of prediction error is chosen as the *best* estimator (see Section 4.2.3).

2. The RKHS meta-model with at most $qmax$ active groups:

The tuning parameter γ is set as zero. A value of μ for which the number of groups in the solution of the RKHS group lasso optimization problem is equal to $qmax$, is computed. This value will be denoted by μ_{qmax} . Then, the RKHS ridge group sparse algorithm is implemented for a grid of values of $\gamma \neq 0$ and μ_{qmax} . This algorithm is described in Section 4.3.2.3.

4.3.1 Calculation of the Gram matrices

The available kernels in the **RKHSMetaMod** package are: linear kernel, quadratic kernel, brownian kernel, matern kernel and gaussian kernel. The usual presentation of these kernels is given in Table 4.1. The choice of the kernel that is done by

Kernel type	Mathematics formula for $u \in \mathbb{R}^n, v \in \mathbb{R}$	RKHSMetaMod name
Linear	$k_a(u, v) = u^T v + 1$	"linear"
Quadratic	$k_a(u, v) = (u^T v + 1)^2$	"quad"
Brownian	$k_a(u, v) = \min(u, v) + 1$	"brownian"
Matern	$k_a(u, v) = (1 + 2 u - v) \exp(-2 u - v)$	"matern"
Gaussian	$k_a(u, v) = \exp(-2\ u - v\ ^2)$	"gaussian"

Table 4.1: List of the reproducing kernels used to construct the RKHS \mathcal{H} .

the user, determines the functional approximation space. For a chosen kernel, the algorithm to calculate the Gram matrices $K_v, v \in \mathcal{P}$ in the **RKHSMetaMod** package, is implemented as `calc_Kv` function. This algorithm is based on three essential points:

- (1) Modify the chosen kernel:

In order to satisfy the conditions of constructing the RKHS \mathcal{H} described in Section 4.2.2, these kernels are modified according to Equation (4.15). Let us take the example of the Brownian kernel:

Example 4.3.1 *The RKHS associated with the brownian kernel $k_a(X_a, X'_a) = \min(X_a, X'_a) + 1$ is well known to be the set,*

$$\mathcal{H}_a = \left\{ f : [0, 1] \rightarrow \mathbb{R} \text{ is absolutely continuous, and } f(0) = 0, \int_0^1 f'(X_a)^2 dX_a < \infty \right\},$$

with the inner product

$$\langle f, h \rangle_{\mathcal{H}_a} = \int_0^1 f'(X_a) h'(X_a) dX_a.$$

The kernel k_{0a} associated with the brownian kernel is calculated as follows,

$$\begin{aligned} k_{0a} &= \min(X_a, X'_a) + 1 - \frac{(\int_0^1 (\min(X_a, U) + 1)dU)(\int_0^1 (\min(X'_a, U) + 1)dU)}{(\int_0^1 \int_0^1 (\min(U, V) + 1)dUdV)}, \\ &= \min(X_a, X'_a) + 1 - \frac{3}{4}(1 + X_a - \frac{X_a^2}{2})(1 + X'_a - \frac{X_a'^2}{2}). \end{aligned}$$

The RKHS associated with the kernel k_{0a} is the set,

$$\mathcal{H}_{0a} = \left\{ f \in \mathcal{H}_a : \int_0^1 f(X_a)dX_a = 0 \right\}.$$

Finally, the RKHS $\mathcal{H} = \mathbb{1} + \sum_{v \in \mathcal{P}} \mathcal{H}_v$ is the following set,

$$\mathcal{H} = \left\{ f : [0, 1]^d \rightarrow \mathbb{R} : f = f_0 + \sum_{v \in \mathcal{P}} f_v(X_v), \text{ with } f_v \in \mathcal{H}_v \right\}.$$

Remark 4.3.1 In this package, the input variables $X = (X_1, \dots, X_d)$ that are uniformly distributed on $[0, 1]^d$ are considered. In order to consider the input variables that are not distributed uniformly, it suffices to modify a part of the function `calc_Kv` related to the calculation of kernels k_{0a} , $a = 1, \dots, d$. For example, for $X = (X_1, \dots, X_d)$ being distributed with law $P_X = \prod_{a=1}^d P_a$ on $\mathcal{X} = \bigotimes_{a=1}^d \mathcal{X}_a \subset \mathbb{R}^d$, the kernel k_{0a} associated with the brownian kernel is calculated as follows,

$$k_{0a} = \min(X_a, X'_a) + 1 - \frac{(\int_{\mathcal{X}_a} (\min(X_a, U) + 1)dP_a)(\int_{\mathcal{X}_a} (\min(X'_a, U) + 1)dP_a)}{(\int_{\mathcal{X}_a} \int_{\mathcal{X}_a} (\min(U, V) + 1)dP_a dP_a)}.$$

The other parts of function `calc_Kv` remain unchanged.

- (2) Calculate the Gram matrices K_v for all v :

First, for all $a = 1, \dots, d$ the Gram matrices K_a associated with kernels k_{0a} are calculated using Equation (4.15),

$$(K_a)_{i,i'} = k_{0a}(X_{ai}, X_{ai'}).$$

Then, for all $v \in \mathcal{P}$, the Gram matrices K_v associated with kernel $k_v = \prod_{a \in v} k_{0a}$ are calculated as follows:

$$K_v = \bigodot_{a \in v} K_a,$$

where \bigodot denotes the Hadamard product.

- (3) Insure the positive definiteness of the matrices K_v :

The output of the function `calc_Kv` is one of the input arguments of the functions associated with the RKHS group lasso and the RKHS ridge group sparse algorithms. As both of these algorithms rely on the positive definiteness

of these matrices, it is mandatory to have $K_v, v \in \mathcal{P}$ that are positive definite. The options, "correction" and "tol", are provided by the function `calc_Kv` in order to insure the positive definiteness of the matrices $K_v, v \in \mathcal{P}$. Let us briefly explain this part of the algorithm:

For each group $v \in \mathcal{P}$, let $\lambda_{v,i}, i = 1, \dots, n$ be the eigenvalues associated with the matrix K_v . Set $\lambda_{v,\max} = \max_i \lambda_{v,i}$ and $\lambda_{v,\min} = \min_i \lambda_{v,i}$. For some fixed value of tolerance "tol", and for each matrix K_v ,

"if $\lambda_{v,\min} < \lambda_{v,\max} \times \text{tol}$ ",

then the "correction" to K_v is done. That is,

"The eigenvalues of K_v are replaced by $\lambda_{v,i} + \text{epsilon}$ ",

where "epsilon" is equal to $\lambda_{v,\max} \times \text{tol}$.

The value of "tol" is set as $1e^{-8}$ by default, but one may consider a smaller or greater value for it depending on the kernel chosen and the value of n .

The function `calc_Kv` is described in Section 4.4.2 and illustrated in Example 4.5.3.

4.3.2 Optimization algorithms

The RKHS meta-model is the solution of one of the optimization problems: the minimization of the RKHS group lasso criterion presented in Equation (4.13) (if $\gamma = 0$), or the minimization of the RKHS ridge group sparse criterion presented in Equation (4.12) (if $\gamma \neq 0$). In the following the algorithms to solve these optimization problems are presented.

4.3.2.1 RKHS group lasso

A popular technique for doing group wise variable selection is group lasso. With this procedure, depending on the value of the tuning parameter μ , an entire group of predictors may drop out of the model. An efficient algorithm for solving group lasso problem is the classical block coordinate descent algorithm (Boyd et al. (2011), Bubeck (2015)). Following the idea of Fu (1998), Yuan and Lin (2006) implemented a block wise descent algorithm for the group lasso penalized least-squares, under the condition that the model matrices in each group are orthonormal. A block coordinate (gradient) descent algorithm for solving the group lasso penalized logistic regression is then developed by Meier et al. (2008). This algorithm is implemented in the `grplasso` R package available from CRAN at <https://cran.r-project.org/web/packages/grplasso/>. Yang and Zou (2015) proposed an unified algorithm, named groupwise majorization descent, for solving the general group lasso learning problems by assuming that the loss function satisfies a quadratic majorization condition. The implementation of their work is done in the `gglasso` R package available at <https://cran.r-project.org/web/packages/gglasso/> from CRAN.

In order to solve the RKHS group lasso optimization problem, the classical block coordinate descent algorithm is used. The minimization of criterion $C_g(f_0, \theta)$ (see Equation (4.13)) is done along each group v at a time. At each step of the algorithm, the criterion $C_g(f_0, \theta)$ is minimized as a function of the current block's parameters, while the parameters values for the other blocks are fixed to their current values. The procedure is repeated until convergence.

This procedure leads to Algorithm 1 (see Section 4.A for more details on this procedure).

In the **RKHSMetaMod** package the Algorithm 1 is implemented as **RKHSgrplasso** function. This function is described in Section 4.4.2 and illustrated in Example 4.5.3.

Algorithm 1 RKHS group lasso algorithm:

```

1: Set  $\theta_0 = [0]_{|\mathcal{P}| \times n}$ 
2: repeat
3:   Calculate  $f_0 = \arg \min_{f_0} C_g(f_0, \theta)$ 
4:   for  $v \in \mathcal{P}$  do
5:     Calculate  $R_v = Y - f_0 - \sum_{v \neq w} K_w \theta_w$ 
6:     if  $\| \frac{2}{\sqrt{n}} K_v^{1/2} R_v \| \leq \mu_g$  then
7:        $\theta_v \leftarrow 0$ 
8:     else
9:        $\theta_v \leftarrow \arg \min_{\theta_v} C_g(f_0, \theta)$ 
10:    end if
11:  end for
12: until convergence

```

4.3.2.2 RKHS ridge group sparse

In order to solve the RKHS ridge group sparse optimization problem, an adapted block coordinate descent algorithm is proposed. This algorithm provides two steps:

Step 1 Initialize the input parameters by the solutions of the RKHS group lasso algorithm for each value of the tuning parameter μ , and implement the RKHS ridge group sparse algorithm through active support of the RKHS group lasso solutions until it achieves convergence.

This step is provided in order to decrease the execution time. In fact, instead of implementing the RKHS ridge group sparse algorithm over the set of all groups \mathcal{P} , it is implemented only over the active support obtained by the RKHS group lasso algorithm, $\widehat{S}_{f_{\text{Group Lasso}}}$.

Step 2 Re-initialize the input parameters with the obtained solutions of Step 1 and implement the RKHS ridge group sparse algorithm through all groups in \mathcal{P} until it achieves convergence.

This second step makes it possible to verify that no group is missing in the output of Step 1.

This procedure leads to Algorithm 2 (see Section 4.A for more details on this procedure).

In the **RKHSMetaMod** package the Algorithm 2 is implemented as `pen_MetMod` function. This function is described in Section 4.4.2 and illustrated in Example 4.5.3.

Algorithm 2 RKHS ridge group sparse algorithm:

- 1: **Step 1:**
 - 2: Set $\theta_0 = \hat{\theta}_{\text{Group Lasso}}$ and $\hat{\mathcal{P}} = \hat{S}_{\hat{f}_{\text{Group Lasso}}}$
 - 3: **repeat**
 - 4: Calculate $f_0 = \arg \min_{f_0} C(f_0, \theta)$
 - 5: **for** $v \in \hat{\mathcal{P}}$ **do**
 - 6: Calculate $R_v = Y - f_0 - \sum_{v \neq w} K_w \theta_w$
 - 7: Solve $J^* = \arg \min_{\hat{t}_v \in \mathbb{R}^n} \{J(\hat{t}_v), \text{ such that } \|K_v^{-1/2} \hat{t}_v\| \leq 1\}$
 - 8: **if** $J^* \leq \gamma$ **then**
 - 9: $\theta_v \leftarrow 0$
 - 10: **else**
 - 11: $\theta_v \leftarrow \arg \min_{\theta_v} C(f_0, \theta)$
 - 12: **end if**
 - 13: **end for**
 - 14: **until** convergence
 - 15: **Step 2:**
 - 16: Implement the same procedure as **Step 1** with $\theta_0 = \hat{\theta}_{\text{old}}$, $\hat{\mathcal{P}} = \mathcal{P} \triangleright \hat{\theta}_{\text{old}}$ is the estimation of θ in **Step 1**.
-

4.3.2.3 RKHS meta-model with *qmax* active groups

By considering some prior information about the data, one may be interested in an RKHS meta-model \hat{f} with the number of active groups not greater than some "*qmax*". In order to obtain the estimator \hat{f} with at most "*qmax*" active groups, the following procedure is provided in the **RKHSMetaMod** package:

- ✓ First, the tuning parameter γ is set as zero and a value of μ for which the solution of the RKHS group lasso algorithm, Algorithm 1, contains exactly *qmax* active groups is computed. This value is denoted by μ_{qmax} .
- ✓ Then, the RKHS ridge group sparse algorithm, Algorithm 2, is implemented by setting the tuning parameter μ equals to μ_{qmax} , and a grid of values of the tuning parameter $\gamma > 0$.

Algorithm 3 Algorithm to estimate RKHS meta-model with at most $qmax$ active groups:

- 1: Calculate $\mu_{\max} = \max_v \frac{2}{\sqrt{n}} \|K_v^{1/2}(Y - \bar{Y})\|$
 - 2: Set $\mu_1 = \mu_{\max}$ and $\mu_2 = \frac{\mu_{\max}}{rat}$ ▷ "rat" is setted by user.
 - 3: **repeat**
 - 4: Implement RKHS group lasso algorithm, Algorithm 1, with $\mu_i = \frac{\mu_1 + \mu_2}{2}$
 - 5: Set $q = |\hat{S}_{\hat{f}_{\text{Group Lasso}}}|$
 - 6: **if** $q > qmax$ **then**
 - 7: Set $\mu_1 = \mu_1$ and $\mu_2 = \mu_i$
 - 8: **else**
 - 9: Set $\mu_1 = \mu_i$ and $\mu_2 = \mu_2$
 - 10: **end if**
 - 11: **until** $q = qmax$ or $i > \text{Num}$ ▷ "Num" is setted by user.
 - 12: Implement RKHS ridge group sparse algorithm, Algorithm 2, with $(\mu = \mu_{qmax}, \gamma > 0)$
-

This procedure leads to Algorithm 3. This algorithm is implemented in the **RKHSMetaMod** package, as function `RKHSMetMod_qmax`. This function is described in Section 4.4.1 and illustrated in Example 4.5.2.

Remark 4.3.2 *As both terms in the penalty function of criterion (4.12) enforce sparsity to the solution, the estimator obtained by solving the RKHS ridge group sparse associated with the pair of the tuning parameters $(\mu_{qmax}, \gamma > 0)$ may contain a smaller number of groups than the solution of the RKHS group lasso optimization problem (i.e. the RKHS ridge group sparse with $(\mu_{qmax}, \gamma = 0)$). And therefore, the estimated RKHS meta-model contains at most "qmax" active groups.*

4.4 Overview of the RKHSMetaMod functions

In the R environment, one can install and load the **RKHSMetaMod** package by using the following commands:

```
R> install.packages("RKHSMetaMod")
R> library("RKHSMetaMod")
```

The optimization problems in this package are solved using block coordinate descent algorithm which requires various computational algorithms including generalized Newton, Broyden and Hybrid methods. In order to gain the efficiency in terms of the calculation time and be able to deal with high dimensional problems, the computationally efficient tools of C++ packages **Eigen** (Guennebaud et al. (2010)) and **GSL** (Galassi (2018)) via **RcppEigen** (Bates and Eddelbuettel (2013)) and **RcppGSL** (Eddelbuettel and Francois (2019)) packages, are used in the **RKHSMetaMod** package. For different examples of usage of **RcppEigen** and **RcppGSL** functions see the

work by [Eddelbuettel \(2013\)](#).

The complete documentation of **RKHSMetaMod** package is available at <https://cran.r-project.org/web/packages/RKHSMetaMod/RKHSMetaMod.pdf>. Here, a brief documentation of some of its main and companion functions is presented in Sections 4.4.1 and 4.4.2, respectively.

4.4.1 Main RKHSMetaMod functions

Let us begin by introducing some notations. For a given $D_{\max} \in \mathbb{N}$, let $\mathcal{P}_{D_{\max}}$ be the set of parts of $\{1, \dots, d\}$ with dimensions 1 to D_{\max} . The cardinal of $\mathcal{P}_{D_{\max}}$ is denoted by v_{\max} ,

$$v_{\max} = \sum_{j=1}^{D_{\max}} \binom{d}{j}.$$

RKHSMetMod function: For a given value of D_{\max} and a chosen kernel (see Table 4.1), this function calculates the Gram matrices K_v , $v \in \mathcal{P}_{D_{\max}}$, and produces a sequence of estimators \hat{f} associated with a given grid of values of tuning parameters μ, γ , i.e. the solutions to the RKHS ridge group sparse (if $\gamma \neq 0$) or the RKHS group lasso problem (if $\gamma = 0$). Table 4.2 gives a summary of all input arguments of the **RKHSMetMod** function and default values for non-mandatory arguments.

The **RKHSMetMod** function returns a list of l components, with l equals to the number of pairs of the tuning parameters (μ, γ) , i.e. $l = |\text{gamma}| \times |\text{frc}|$. Each component of the list is a list of three components "mu", "gamma" and "Meta-Model":

- ✓ mu: value of the tuning parameter μ if $\gamma > 0$, or $\mu_g = \sqrt{n} \times \mu$ if $\gamma = 0$.
- ✓ gamma: value of the tuning parameter γ .
- ✓ Meta-Model: an RKHS ridge group sparse or RKHS group lasso object associated with the tuning parameters mu and gamma. Table 4.3 gives a summary of all arguments of the output "Meta-Model" of **RKHSMetMod** function.

RKHSMetMod_qmax function: For a given value of D_{\max} and a chosen kernel (see Table 4.1), this function calculates the Gram matrices K_v , $v \in \mathcal{P}_{D_{\max}}$, determines μ , denoted μ_{qmax} , for which the number of active groups in the RKHS group lasso solution is equal to $qmax$, and produces a sequence of estimators \hat{f} associated with the tuning parameter μ_{qmax} and a grid of values of the tuning parameter γ . All the estimators \hat{f} produced by this function have at most $qmax$ active groups in their support. This function has the following input arguments:

- Y, X , kernel, D_{\max} , gamma, verbose (see Table 4.2).
- qmax: integer, the maximum number of active groups in the obtained solution.

Input parameter	Description
Y	Vector of the response observations of size n .
X	Matrix of the input observations with n rows and d columns. Rows correspond to the observations and columns correspond to the variables.
kernel	Character, indicates the type of the kernel (see Table 4.1) chosen to construct the RKHS \mathcal{H} .
Dmax	Integer, between 1 and d , indicates the maximum order of interactions considered in the RKHS meta-model: Dmax= 1 is used to consider only the main effects, Dmax= 2 to include the main effects and the second-order interactions, and so on.
gamma	Vector of non-negative scalars, values of the tuning parameter γ in decreasing order. If $\gamma = 0$ the function solves the RKHS group lasso optimization problem and for $\gamma > 0$ it solves the RKHS ridge group sparse optimization problem.
frc	Vector of positive scalars. Each element of the vector sets a value to the tuning parameter μ : $\mu = \mu_{\max}/(\sqrt{n} \times \text{frc})$. The value μ_{\max} (see Equation (4.16)) is calculated inside the program.
verbose	Logical. Set as TRUE to print: the group v for which the correction of the Gram matrix K_v is done (see Section 4.3.1), and for each pair of the tuning parameters (μ, γ) : the number of current iteration, active groups and convergence criterion. It is set as FALSE by default.

Table 4.2: List of the input arguments of the RKHSMetMod function.

- rat: positive scalar, to restrict the minimum value of μ considered in Algorithm 3,

$$\mu_{\min} = \frac{\mu_{\max}}{(\sqrt{n} \times \text{rat})},$$

where the value of μ_{\max} is given by Equation (4.16) and is calculated inside the program.

- Num: integer, to restrict the number of different values of the tuning parameter μ to be evaluated in the RKHS group lasso algorithm until it achieves μ_{qmax} . For example, if Num equals to 1 the program is implemented for three different

Output parameter	Description
intercept	Scalar, estimated value of intercept.
teta	Matrix with vMax rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, vMax$.
fit.v	Matrix with n rows and vMax columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
fitted	Vector of size n , indicates the estimator of m .
Norm.n	Vector of size vMax, estimated values for the empirical L^2 -norm.
Norm.H	Vector of size vMax, estimated values for the Hilbert norm.
supp	Vector of active groups.
Nsupp	Vector of the names of the active groups.
SCR	Scalar equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
crit	Scalar indicates the value of the penalized criterion.
gamma.v	Vector of size vMax, coefficients of the empirical L^2 -norm.
mu.v	Vector of size vMax, coefficients of the Hilbert norm.
iter	List of two components: maxIter, and the number of iterations until the convergence is achieved.
convergence	TRUE or FALSE. Indicates whether the algorithm has converged or not.
RelDiffCrit	Scalar, value of the first convergence criterion at the last iteration, i.e. $\ \frac{\theta_{lastIter} - \theta_{lastIter-1}}{\theta_{lastIter-1}}\ ^2$.
RelDiffPar	Scalar, value of the second convergence criterion at the last iteration, i.e. $\frac{crit_{lastIter} - crit_{lastIter-1}}{crit_{lastIter-1}}$.

Table 4.3: List of the arguments of the output "Meta-Model" of RKHSMetMod function.

values of $\mu \in [\mu_{\min}, \mu_{\max}]$:

$$\begin{aligned} \mu_1 &= \frac{(\mu_{\min} + \mu_{\max})}{2}, \\ \mu_2 &= \begin{cases} \frac{(\mu_{\min} + \mu_1)}{2} & \text{if } |\widehat{S}_{\widehat{f}(\mu_1)_{\text{Group Lasso}}}| < qmax, \\ \frac{(\mu_1 + \mu_{\max})}{2} & \text{if } |\widehat{S}_{\widehat{f}(\mu_1)_{\text{Group Lasso}}}| > qmax, \end{cases} \\ \mu_3 &= \mu_{\min}, \end{aligned}$$

where $|\widehat{S}_{\widehat{f}(\mu_1)_{\text{Group Lasso}}}|$ is the number of active groups in the solution of the RKHS group lasso problem, Algorithm 1, associated with μ_1 .

If $\text{Num} > 1$, the path to cover the interval $[\mu_{\min}, \mu_{\max}]$ is detailed in Algorithm 3.

The RKHSMetMod_qmax function returns a list of three components "mus", "qs", and "MetaModel":

- ✓ mus: vector of all values of μ_i in Algorithm 3.
- ✓ qs: vector with the same length as mus. Each element of the vector shows the number of active groups in the RKHS meta-model obtained by solving RKHS group lasso problem for an element in mus.
- ✓ MetaModel: list with the same length as the vector gamma. Each component of the list is a list of three components "mu", "gamma" and "Meta-Model":
 - mu: value of μ_{qmax} .
 - gamma: element of the input vector gamma associated with the estimated "Meta-Model".
 - Meta-Model: an RKHS ridge group sparse or RKHS group lasso object associated with the tuning parameters mu and gamma (see Table 4.3).

4.4.2 Companion functions

calc_Kv function: For a given value of Dmax and a chosen kernel (see Table 4.1), this function calculates the Gram matrices K_v , $v \in \mathcal{P}_{Dmax}$, and returns their associated eigenvalues and eigenvectors. This function has,

- ✓ four mandatory input arguments:
 - Y, X, kernel, Dmax (see Table 4.2).
- ✓ three facultative input arguments:
 - correction: logical, set as TRUE to make correction to the matrices K_v (see Section 4.3.1). It is set as TRUE by default.
 - verbose: logical, set as TRUE to print: the group for which the correction is done. It is set as TRUE by default.
 - tol: scalar to be chosen small, set as $1e^{-8}$ by default.

The calc_Kv function returns a list of two components "kv" and "names.Grp":

- ✓ kv: list of vMax components, each component is a list of,
 - Evalues: vector of eigenvalues.
 - Q: matrix of eigenvectors.
- ✓ names.Grp: vector of group names of size vMax.

RKHSgrplasso function: For a given value of the tuning parameter μ_g , this function fits the solution to the RKHS group lasso optimization problem by implementing Algorithm 1. This function has,

- ✓ three mandatory input arguments:

- Y (see Table 4.2).
 - Kv : list of the eigenvalues and the eigenvectors of the positive definite Gram matrices K_v for $v = 1, \dots, vMax$ and their associated group names (output of the function `calc_Kv`).
 - μ : positive scalar indicates the value of the tuning parameter μ_g defined in Equation (4.14).
- ✓ two facultative input arguments:
- `maxIter`: integer, to set the maximum number of loops through all groups. It is set as 1000 by default.
 - `verbose`: logical, set as TRUE to print: the number of current iteration, active groups and convergence criterion. It is set as FALSE by default.

This function returns an RKHS group lasso object associated with the tuning parameter μ_g . Its output is a list of 13 components:

- ✓ `intercept`, `teta`, `fit.v`, `fitted`, `Norm.H`, `supp`, `Nsupp`, `SCR`, `crit`, `MaxIter`, `convergence`, `RelDiffCrit`, `RelDiffPar` (see Table 4.3).

mu_max function: This function calculates the value μ_{max} defined in Equation (4.16). It has two mandatory input arguments: the response vector Y , and the list `matZ` of the eigenvalues and eigenvectors of the positive definite Gram matrices K_v for $v = 1, \dots, vMax$. This function returns the μ_{max} value.

pen_MetMod function: This function produces a sequence of the RKHS meta-models associated with a given grid of values of the tuning parameters μ, γ . Each RKHS meta-model in the sequence is the solution to the RKHS ridge group sparse optimization problem (obtained by implementing Algorithm 2) associated with a pair of values of (μ, γ) in the grid of values of μ, γ . This function has,

- ✓ seven mandatory input arguments:
- Y (see Table 4.2).
 - `gamma`: vector of positive scalars. Values of the penalty parameter γ in decreasing order.
 - Kv : list of the eigenvalues and the eigenvectors of the positive definite Gram matrices K_v for $v = 1, \dots, vMax$ and their associated group names (output of the function `calc_Kv`).
 - μ : vector of positive scalars. Values of the tuning parameter μ in decreasing order.
 - `resg`: list of the RKHS group lasso objects associated with the components of "mu", used as initial parameters at Step 1.

- `gama_v` and `mu_v`: vector of `vMax` positive scalars. These two inputs are optional, they are provided to associate the weights to the two penalty terms in the RKHS ridge group sparse criterion (4.12). They set to scalar 0, to consider no weights, i.e. all weights equal to 1.

✓ three facultative input arguments:

- `maxIter`: integer, to set the maximum number of loops through initial active groups at Step 1 and maximum number of loops through all groups at Step 2. It is set as 1000 by default.
- `verbose`: logical, set as TRUE to print: for each pair of the tuning parameters (μ, γ) : the number of current iteration, active groups and convergence criterion. It is set as FALSE by default.
- `calcStwo`: logical, set as TRUE to execute Step 2. It is set as FALSE by default.

The function `pen_MetMod` returns a list of l components, with l equals to the number of pairs of the tuning parameters (μ, γ) . Each component of the list is a list of three components "mu", "gamma" and "Meta-Model":

- ✓ `mu`: positive scalar, an element of the input vector "mu" associated with the estimated "Meta-Model".
- ✓ `gamma`: positive scalar, an element of the input vector "gamma" associated with the estimated "Meta-Model".
- ✓ `Meta-Model`: an RKHS ridge group sparse object associated with the tuning parameters `mu` and `gamma` (see Table 4.3).

PredErr function: By considering a testing dataset, this function calculates the prediction errors for the obtained RKHS meta-models. This function has eight mandatory input arguments:

- `X`, `gamma`, `kernel`, `Dmax` (see Table 4.2).
- `XT`: matrix of observations of the testing dataset with n^{test} rows and d columns.
- `YT`: vector of response observations of the testing dataset of size n^{test} .
- `mu`: vector of positive scalars. Values of the tuning parameter μ in decreasing order.
- `res`: list of the estimated RKHS meta-models for the learning dataset associated with the tuning parameters (μ, γ) (it could be the output of one of the functions `RKHSMetMod`, `RKHSMetMod_qmax` or `pen_MetMod`).

Note that, the same kernel and Dmax have to be chosen as the ones used for the learning dataset.

The function `PredErr` returns a matrix of the prediction errors. Each element of the matrix corresponds to the prediction error of one RKHS meta-model in "res".

SI_emp function: For each RKHS meta-model \hat{f} , this function calculates the empirical Sobol indices for all groups that are active in the support of \hat{f} . This function has two input arguments:

- res: list of the estimated meta-models using RKHS ridge group sparse or RKHS group lasso algorithms (it could be the output of one of the functions `RKHSMetMod`, `RKHSMetMod_qmax` or `pen_MetMod`).
- ErrPred: matrix or NULL. If matrix, each element of the matrix corresponds to the prediction error of an RKHS meta-model in "res" (output of the function `PredErr`). Set as NULL by default.

The empirical Sobol indices are then calculated for each RKHS meta-model in "res", and a list of vectors of the Sobol indices is returned.

If the argument "ErrPred" is the matrix of the prediction errors, the vector of empirical Sobol indices is returned for the *best* RKHS meta-model in the "res".

4.5 RKHSMetaMod through examples

Let us consider the g-function of Sobol (Saltelli et al. (2009)) in the Gaussian regression framework, i.e.

$$Y = m(X) + \sigma\varepsilon, \sigma > 0,$$

where the error term ε is a centered Gaussian random variable, and where the function m is the g-function of Sobol defined over $[0, 1]^d$ by,

$$m(X) = \prod_{a=1}^d \frac{|4x_a - 2| + c_a}{1 + c_a}, c_a > 0. \quad (4.17)$$

The Sobol indices of the g-function can be expressed analytically:

$$\forall v \in \mathcal{P}, S_v = \frac{1}{D} \prod_{a \in v} D_a, D_a = \frac{1}{3(1 + c_a)^2}, D = \prod_{a=1}^d (D_a + 1) - 1.$$

Set $c_1 = 0.2$, $c_2 = 0.6$, $c_3 = 0.8$ and $(c_a)_{a>3} = 100$. With these values of coefficients c_a , the variables X_1, X_2 and X_3 explain 99.99% of the variance of the function $m(X)$ (Durrande et al. (2013)). The values of S_v , $v \in \mathcal{P}$, when $d = 5$ and $d = 10$ are displayed in Tables 4.6 and 4.7, respectively.

In this Section, four examples are presented. In all examples the value of Dmax is set as three. Example 4.5.1 illustrates the use of the `RKHSMetMod` function by considering three different kernels, "matern", "brownian", and "gaussian" (see Table

4.1), and three datasets of $n \in \{50, 100, 200\}$ observations and $d = 5$ input variables. In Example 4.5.2, the function `RKHSMetMod_qmax` is illustrated for dataset of $n = 500$ observations and $d = 10$ input variables. The larger datasets with $n \in \{1000, 2000, 5000\}$ observations and $d = 10$ input variables are studied in Examples 4.5.3 and 4.5.4.

In each example, two independent datasets are generated: (X, Y) to estimate the meta-models, and (XT, YT) to estimate the prediction errors. The design matrices X and XT are the Latin Hypercube Samples of the input variables that are generated using `maximinLHS` function of the package `lhs` available at <https://CRAN.R-project.org/package=lhs>:

```
R> library(lhs)
R> X <- maximinLHS(n,d)
R> XT <- maximinLHS(n,d)
The response variables  $Y$  and  $YT$  are calculated as  $Y = m(X) + \sigma\varepsilon$  and  $YT = m(XT) + \sigma\varepsilon_T$ , where  $\sigma = 0.2$ , and  $\varepsilon, \varepsilon_T$  are distributed independently according to the centered Gaussian distribution with variance equals to one:
R> a <- c(0.2,0.6,0.8,100,100,100,100,100,100,100)[1:d]
R> sigma <- 0.2
R> g=1;for (i in 1:d) g=g*(abs(4*X[,i]-2)+a[i])/(1+a[i])
R> epsilon <- rnorm(n,0,1)
R> Y <- g + sigma*epsilon
R> gT=1;for (i in 1:d) gT=gT*(abs(4*XT[,i]-2)+a[i])/(1+a[i])
R> epsilonT <- rnorm(n,0,1)
R> YT <- gT + sigma*epsilonT
```

Example 4.5.1 *RKHS meta-model estimation using `RKHSMetMod` function:*

In this example, three datasets of n points `maximinLHS` over $[0, 1]^d$ are generated with $n \in \{50, 100, 200\}$ and $d = 5$, and a grid of five values for each of the tuning parameters μ and γ is considered as follows:

$$\mu_{(1:5)} = \frac{\mu_{max}}{(\sqrt{n} \times 2^{(2:6)})}, \gamma_{(1:5)} = (0.2, 0.1, 0.01, 0.005, 0).$$

For each dataset, the experiment is repeated $N_r = 50$ times. At each repetition, the RKHS meta-models associated with the pair of the tuning parameters (μ, γ) are estimated using the `RKHSMetMod` function:

```
R> kernel <- "matern" # kernel <- "brownian" # kernel <- "gaussian"
R> Dmax <- 3
R> gamma <- c(0.2,0.1,0.01,0.005,0)
R> frc <- c(4,8,16,32,64)
R> res <- RKHSMetMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
```

These meta-models are evaluated using a testing dataset. The prediction errors are computed for them using the `PredErr` function. The RKHS meta-model with minimum prediction error is chosen to be the *best* estimator for the model. Finally,

the Sobol indices are computed for the *best* RKHS meta-model using the function `SI_emp`:

```
R> l <- length(gamma)
R> mu <- vector();for(i in 1:length(frc)){mu[i] <- res[[(i-1)*l+1]]$mu}
R> Err <- PredErr(X,XT,YT,mu,gamma,res,kernel,Dmax)
R> SI <- SI_emp(res,Err)
```

The performances of this method for estimating a meta-model are evaluated by considering a third dataset $(m(X_i^{third}), X_i^{third})$, $i = 1, \dots, N$, with $N = 1000$. The global prediction error is calculated as follows:

Let $\hat{f}_r(\cdot)$ be the *best* RKHS meta-model obtained in the repetition r , $r = 1, \dots, N_r$, then

$$GPE = \frac{1}{N_r} \sum_{r=1}^{N_r} \frac{1}{N} \sum_{i=1}^N (\hat{f}_r(X_i^{third}) - m(X_i^{third}))^2.$$

The values of GPE obtained for different kernels and values of n are given in Table 4.4. As expected the value of GPE decreases as n increases. The lowest values of

n	50	100	200
GPE_m	0.13	0.07	0.03
GPE_b	0.14	0.10	0.05
GPE_g	0.15	0.10	0.07

Table 4.4: Example 4.5.1: The columns of the table correspond to the different datasets with $n \in \{50, 100, 200\}$ and $d = 5$. Each line of the table, from up to down, gives the value of GPE obtained for each dataset associated with the "matern", "brownian" and "gaussian" kernels, respectively.

GPE are obtained when using the "matern" kernel.

In order to sum up the behaviour of the procedure for estimating the Sobol indices, the mean square error (MSE) is estimated as follows:

Let

$$b_v^2 = (\hat{S}_{v,\cdot} - S_v)^2, \text{ and } w_v^2 = \frac{1}{N_r} \sum_{r=1}^{N_r} (\hat{S}_{v,r} - \hat{S}_{v,\cdot})^2,$$

where for each group v , S_v denotes the true values of the Sobol indices, and for $\hat{S}_{v,r}$ being the empirical Sobol indices of the *best* RKHS meta-model in repetition r , $\hat{S}_{v,\cdot}$ denotes the mean of the empirical Sobol indices of the *best* RKHS meta-models through all repetitions:

$$\hat{S}_{v,\cdot} = \frac{1}{N_r} \sum_{r=1}^{N_r} \hat{S}_{v,r}.$$

Then,

$$MSE = \sum_v (b_v^2 + w_v^2).$$

**Chapter 4. Estimate the Hoeffding decomposition of a complex model
140 by solving RKHS ridge group sparse optimization problem**

The obtained values of MSE for different kernels and values of n are given in Table 4.5. As expected, the values of MSE are smaller for larger values of n . The smallest

n	50	100	200
MSE_m	75.12	46.72	28.22
MSE_b	110.71	84.99	41.06
MSE_g	78.22	94.67	67.02

Table 4.5: Example 4.5.1: The columns of the table correspond to the different datasets with $n \in \{50, 100, 200\}$ and $d = 5$. Each line of the table, from up to down, gives the value of MSE obtained for each dataset associated with the "matern", "brownian" and "gaussian" kernels, respectively.

values are obtained when using "matern" kernel.

The means of the empirical Sobol indices of the *best* RKHS meta-models through all repetitions for $n = 200$ and "matern" kernel are displayed in Table 4.6. It

v	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}	sum
S_v	43.24	24.32	19.22	5.63	4.45	2.50	0.58	99.94
$\widehat{S}_{v..}$	46.10	26.33	20.62	2.99	2.22	1.13	0.0	99.39

Table 4.6: Example 4.5.1: The first line of the table gives the true values of the Sobol indices $\times 100$. The second line gives the mean of the estimated empirical Sobol indices $\times 100$ greater than 10^{-2} calculated over fifty simulations for $n = 200$ and "matern" kernel. The sum of the Sobol indices is displayed in the last column.

appears that the estimated Sobol indices are close to the true ones, nevertheless they are over estimated for the main effects, i.e. groups $v \in \{\{1\}, \{2\}, \{3\}\}$, and under estimated for the interactions of order two and three, i.e. groups $v \in \{\{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$.

Note that, the strategy of choosing the tuning parameters is based on the minimization of the prediction error of the estimated meta-model, which may not minimize the error of estimating the Sobol indices.

Taking into account the results obtained for this Example 4.5.1, the calculations in the rest of the examples is done using only the "matern" kernel.

Example 4.5.2 Estimate the meta-models with at most "qmax" active groups:

A dataset of n points maximLHS over $[0, 1]^d$ with $n = 500$ and $d = 10$ is generated. The true values of the Sobol indices with $d = 10$ are displayed in Table 4.7. As we

v	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}	sum
S_v	43.26	24.33	19.22	5.63	4.45	2.50	0.58	99.97

Table 4.7: Example 4.5.2: The true values of the Sobol indices $\times 100$ when $d = 10$.

can see, the main factors X_1, X_2 , and X_3 explain almost all of the variability in the

model. So, one may be interested in estimating the function $m(X)$ (see Equation (4.17)) by a meta-model that includes at most three active groups (the main effects only). In order to calculate the RKHS meta-models that contain at most three active groups the `RKHSMetMod_qmax` function is used with,

✓ "gamma" = (0.2, 0.1, 0.01, 0.005, 0),

✓ "rat" = 100: the minimum value of μ considered in the algorithm is then

$$\mu_{min} = \frac{\mu_{max}}{(\sqrt{n} \times 100)},$$

✓ "Num" = 10: the maximum number of values of $\mu \in [\mu_{min}, \mu_{max})$ to be evaluated is equal to twelve (see Algorithm 3).

```
R> kernel <- "matern"
R> Dmax <- 3
R> gamma <- c(0.2,0.1,0.01,0.005,0)
R> qmax <- 3;Num <- 10;rat <- 100
R> res <- RKHSMetMod_qmax(Y,X,kernel,Dmax,gamma,qmax,Num,rat,FALSE)
The RKHS meta-models are estimated for the obtained value of  $\mu_{qmax}$  and different
values of the tuning parameter  $\gamma$ :
R> for(i in 1:length(gamma)){
print(paste("In meta model ",i))
print(paste("the value of mu is: ",res$MetaModel[[i]]$mu,
"and the value of gamma is: ",res$MetaModel[[i]]$gamma))
print("the active groups are: ")
print(res$MetaModel[[i]]$'Meta-Model'$Nsupp)
}
"In meta model 1"
"the value of mu is: 0.093 and the value of gamma is: 0.2"
"the active groups are: "
"v1." "v2." "v3."
"In meta model 2"
"the value of mu is: 0.093 and the value of gamma is: 0.1"
"the active groups are: "
"v1." "v2." "v3."
"In meta model 3"
"the value of mu is: 0.093 and the value of gamma is: 0.01"
"the active groups are: "
"v1." "v2." "v3."
"In meta model 4"
"the value of mu is: 0.093 and the value of gamma is: 0.005"
"the active groups are: "
"v1." "v2." "v3."
"In meta model 5"
```



```
"the value of mu is: 2.083 and the value of gamma is: 0"
"the active groups are: "
"v1." "v2." "v3."
```

Let us comment the outputs of the function `RKHSMetMod_qmax`: for $\gamma \neq 0$ the value "mu" corresponds to the value of $\mu_{qmax=3}$ which is equal to 0.093, while for $\gamma = 0$ the value "mu" corresponds to the value of μ_g defined in Equation (4.14),

$$\mu_g = \sqrt{n} \times 0.093 = 2.083.$$

For each pair of the tuning parameters (μ_{qmax}, γ_i) , $i = 1, \dots, 5$, the estimated RKHS meta-model contains three groups. The groups associated with X_1 , X_2 , and X_3 are "v1.", "v2.", and "v3.", that are active in the estimators obtained, as expected.

Example 4.5.3 *A time saving trick to obtain the "optimal" tuning parameters when dealing with larger datasets:*

A dataset of n points `maximinLHS` over $[0, 1]^d$ with $n = 1000$ and $d = 10$ is generated. Firstly, the eigenvalues and eigenvectors of the positive definite matrices K_v , and the value of μ_{max} is computed using functions `calc_Kv` and `mu_max`, respectively:

```
R> kernel <- "matern"
R> Dmax <- 3
R> Kv <- calc_Kv(X, kernel, Dmax, TRUE, TRUE)
R> mumax <- mu_max(Y, Kv$kv)
```

Then, the two following steps are considered:

1. Set $\gamma = 0$ and,

$$\mu_{(1:9)} = \frac{\mu_{max}}{(\sqrt{n} \times 2^{(2:10)})}.$$

Calculate the RKHS meta-models associated with the values of $\mu_g = \mu \times \sqrt{n}$ by using the function `RKHSgrplasso`. Gather the obtained RKHS meta-models in a list, "res_g". While this job could be done with the function `RKHSMetMod` by setting $\gamma = 0$, in this example we use the function `RKHSgrplasso` in order to avoid the re-calculation of K_v 's at the next step. Thereafter, for each estimator in the `res_g` the prediction error is calculated by considering a new dataset and using the function `PredErr`. The value of μ with the smallest error of prediction in this step is denoted by μ_i .

Let us implement this step:

For a grid of values of μ_g , a sequence of the RKHS meta-models are calculated and gathered in the "res_g" list:

```
R> frc <- c(4,8,16,32,64,128,256,512,1024)
R> mu_g <- mumax/frc
R> res_g <- list();resg <- list()
R> for(i in 1:length(mu_g)){
resg[[i]] <- RKHSgrplasso(Y,Kv,mu_g[i],1000,FALSE)
res_g[[i]] <- list("mu_g"=mu_g,"gamma"=0,"MetaModel"=resg[[i]])
```

```
}

```

Output `res_g` contains nine RKHS meta-models and they are evaluated using a testing dataset:

```
R> gamma <- c(0)

```

```
R> Err_g <- PredErr(X,XT,YT,mu_g,gamma,res_g,kernel,Dmax)

```

The prediction errors of the RKHS meta-models obtained in this step are displayed in Table 4.8. It appears that the minimum prediction error corre-

μ_g	1.304	0.652	0.326	0.163	0.081	0.041	0.020	0.010	0.005
$\gamma = 0$	0.197	0.156	0.145	0.097	0.063	0.055	0.056	0.063	0.073

Table 4.8: Example 4.5.3: Obtained prediction errors in step 1.

sponds to the solution of the RKHS group lasso algorithm with $\mu_g = 0.041$, so $\mu_i = 0.041/\sqrt{n}$.

- Choose a smaller grid of values of μ , $(\mu_{(i-1)}, \mu_i, \mu_{(i+1)})$, and set a grid of values of $\gamma > 0$. Calculate the RKHS meta-models associated with each pair of the tuning parameters (μ, γ) by the function `pen_MetMod`. Calculate the prediction errors for the new sequence of the RKHS meta-models using the function `PredErr`. Compute the empirical Sobol indices for the *best* estimator.

Let us go back to the implementation of the example and apply this step 2:

The grid of values of μ in this step is,

$$\left(\frac{0.081}{\sqrt{n}}, \frac{0.041}{\sqrt{n}}, \frac{0.020}{\sqrt{n}}\right).$$

The RKHS meta-models associated with this grid of values of μ are gathered in a new list "resgnew". Set $\gamma_{(1:4)} = (0.2, 0.1, 0.01, 0.005)$, and calculate the RKHS meta-models for this new grid of values of (μ, γ) using `pen_MetMod` function:

```
R> gamma <- c(0.2,0.1,0.01,0.005)

```

```
R> mu <- c(mu_g[5],mu_g[6],mu_g[7])/sqrt(n)

```

```
R> resgnew <- list()

```

```
R> resgnew[[1]] <- resg[[5]];resgnew[[2]] <- resg[[6]];resgnew[[3]]
<- resg[[7]]

```

```
R> res <- pen_MetMod(Y,Kv,gamma,mu,resgnew,0,0)

```

The output "res" is a list of twelve RKHS meta-models. These meta-models are evaluated using a new dataset, and their prediction errors are displayed in Table 4.9.

The minimum prediction error is associated with the pair $(0.020/\sqrt{n}, 0.01)$, and the *best* RKHS meta-model is then $\hat{f}_{(0.020/\sqrt{n}, 0.01)}$.

The performances of this procedure for estimating the Sobol indices is evaluated using the relative error (RE) defined as follows:

μ	$0.081/\sqrt{n}$	$0.041/\sqrt{n}$	$0.020/\sqrt{n}$
$\gamma = 0.2$	0.153	0.131	0.119
$\gamma = 0.1$	0.098	0.079	0.072
$\gamma = 0.01$	0.065	0.054	0.053
$\gamma = 0.005$	0.064	0.054	0.054

Table 4.9: Example 4.5.3: Obtained prediction errors in step 2.

For each v , let S_v be the true value of the Sobol indices displayed in Table 4.7 and \widehat{S}_v be the estimated empirical Sobol indices. Then

$$RE = \sum_v \frac{|\widehat{S}_v - S_v|}{S_v}. \quad (4.18)$$

In Table 4.10 the estimated empirical Sobol indices, their sum, and the value of RE are displayed.

v	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}	sum	RE
\widehat{S}_v	42.91	25.50	20.81	4.40	3.84	2.13	0.00	99.60	1.64

Table 4.10: Example 4.5.3: The estimated empirical Sobol indices $\times 100$ greater than 10^{-2} . The last two columns show $\sum_v \widehat{S}_v$ and RE, respectively.

The RE for each group v is smaller than 1.64%, so the estimated Sobol indices in this example are very close to the true values of the Sobol indices displayed in the Table 4.7. In this example the significant values of the Sobol indices for interactions of order two are obtained.

Example 4.5.4 *Dealing with larger datasets:*

Two datasets of n points `maximinLHS` over $[0, 1]^d$ with $n \in \{2000, 5000\}$ and $d = 10$ are generated. In order to obtain one RKHS meta-model associated with one pair of the tuning parameters (μ, γ) , the number of coefficients to be estimated is equal to $n \times v_{\text{Max}} = n \times 175$. Table 4.11 gives the execution time for different functions used throughout the Examples 4.5.1-4.5.4. As we can see, the execution time increases fastly as n increases. In Figure 4.1 the plot of the logarithm of the time versus the logarithm of n is displayed for the functions `calc_Kv`, `mu_max`, `RKHSgrplasso` and `pen_MetMod`. It appears that, the algorithms of these functions are of polynomial time $O(n^\alpha)$ with $\alpha \simeq 3$ for the functions `calc_Kv` and `mu_max`, and $\alpha \simeq 2$ for the functions `RKHSgrplasso` and `pen_MetMod`.

Taking into account the results obtained for the prediction error and the values of $(\widehat{\mu}, \widehat{\gamma})$ in Example 4.5.3, in this example only two values of the tuning parameter μ and one value of the tuning parameter γ are considered:

$$\mu = \left(\frac{\mu_{\max}}{(\sqrt{n} \times 2^7)}, \frac{\mu_{\max}}{(\sqrt{n} \times 2^8)} \right) \text{ and } \gamma = 0.01.$$

(n, d)	calc_Kv	mu_max	RKHSgrplasso	pen_MetMod	$ S_{\hat{f}} $	sum
(100,5)	0.09s	0.01s	1s	2s	18	$\sim 3s$
			2s	3s	19	$\sim 5s$
(500,10)	33s	9s	247s	333s	39	$\sim 10min$
			599s	816s	64	$\sim 24min$
(1000,10)	197s	53s	959s	1336s	24	$\sim 42min$
			2757s	4345s	69	$\sim 2h$
(2000,10)	1498s	420s	3984s	4664s	12	$\sim 2h:56min$
			12951s	22385s	30	$\sim 10h:20min$
(5000,10)	34282s	6684s	38957s	49987s	11	$\sim 36h:05min$
			99221s	111376s	15	$\sim 69h:52min$

Table 4.11: Example 4.5.4: The kernel used is "matern". The execution time for the functions `RKHSgrplasso` and `pen_MetMod` is displayed in each row for two pair of values of tuning parameters ($\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7), \gamma = 0.01$) on up, and ($\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8), \gamma = 0.01$) on below. In the column $|S_{\hat{f}}|$, the number of the active groups associated with each estimated RKHS meta-model is displayed.

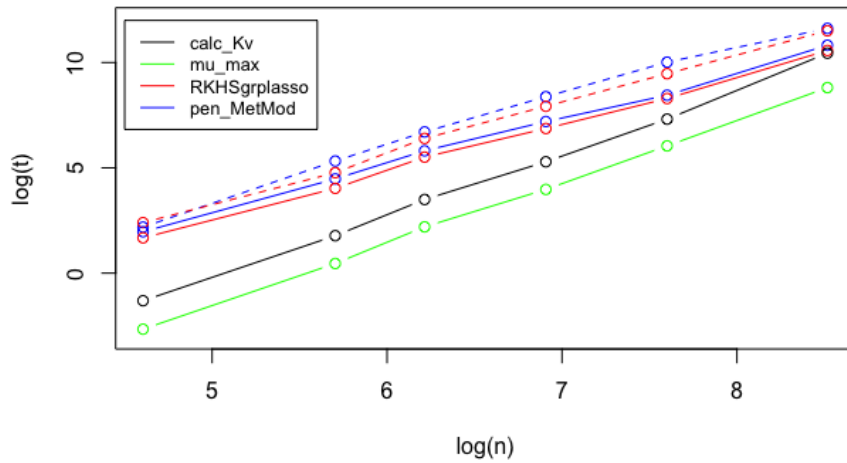


Figure 4.1: Example 4.5.4: Timing plot for $d = 10$, $n \in \{100, 300, 500, 1000, 2000, 5000\}$, and different functions of the `RKHSMetaMod` package. The execution time for the functions `RKHSgrplasso` and `pen_MetMod` is displayed for two pair of values of tuning parameters ($\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7), \gamma = 0.01$) in solid lines, and ($\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8), \gamma = 0.01$) in dashed lines.

The RKHS meta-models associated with the pair of values (μ_i, γ) , $i = 1, 2$ are estimated using the `RKHSMetaMod` function:

**Chapter 4. Estimate the Hoeffding decomposition of a complex model
146 by solving RKHS ridge group sparse optimization problem**

```
R> kernel <- "matern"
R> Dmax <- 3
R> gamma <- c(0.01)
R> frc <- c(128,256)
R> res <- RKHSMetMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
The prediction error and the empirical Sobol indices are then calculated for the
obtained meta-models using the functions PredErr and SI_emp:
R> mu <- vector();mu[1] <- res[[1]]$mu;mu[2] <- res[[2]]$mu
R> Err <- PredErr(X,XT, YT,mu,gamma, res, kernel,Dmax)
R> SI <- SI_emp(res, NULL)
```

The result of the prediction errors associated with the obtained estimators for two different values of n are displayed in Table 4.12. For n equals to 5000 we got smaller

n	$(\mu_{max}/(\sqrt{n} \times 2^7), \gamma)$	$(\mu_{max}/(\sqrt{n} \times 2^8), \gamma)$
2000	0.052	0.049
5000	0.049	0.047

Table 4.12: Example 4.5.4: Obtained prediction errors.

values of the prediction error, so as expected, the prediction quality improves by increasing the number of the observations n . Table 4.13 gives the estimated empirical Sobol indices as well as their sum and the values of RE (see Equation (4.18)). Comparing the values of RE, we can see that the empirical Sobol indices are better

n	v	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}	sum	RE
2000	$\widehat{S}_{v;(\mu_1,\gamma)}$	45.54	24.78	21.01	3.96	3.03	1.65	0.00	99.97	2.12
	$\widehat{S}_{v;(\mu_2,\gamma)}$	45.38	25.07	19.69	4.36	3.66	1.79	0.00	99.95	1.79
5000	$\widehat{S}_{v;(\mu_1,\gamma)}$	44.77	25.39	20.05	4.49	3.38	1.90	0.00	99.98	1.81
	$\widehat{S}_{v;(\mu_2,\gamma)}$	43.78	24.99	19.56	5.43	3.90	2.32	0.00	99.98	1.29

Table 4.13: Example 4.5.4: The estimated empirical Sobol indices $\times 100$ greater than 10^{-2} associated with each estimated RKHS meta-model is printed. The last two columns show $\sum_v \widehat{S}_v$ and RE, respectively. We have $\mu_1 = \mu_{max}/(\sqrt{n} \times 2^7)$, $\mu_2 = \mu_{max}/(\sqrt{n} \times 2^8)$ and $\gamma = 0.01$.

estimated for n equals to 5000, so as expected, the estimation of the Sobol indices is better for larger values of n .

In Figure 4.2 the result of the prediction quality and the Sobol indices for dataset with n equals to 5000, d equals to 10, and (μ_2, γ) are displayed. The line $y = x$ in red crosses the cloud of points as long as the values of the g-function are smaller than three. When the values of the g-function are greater than three, the estimator \widehat{f} tends to under estimate the g-function.

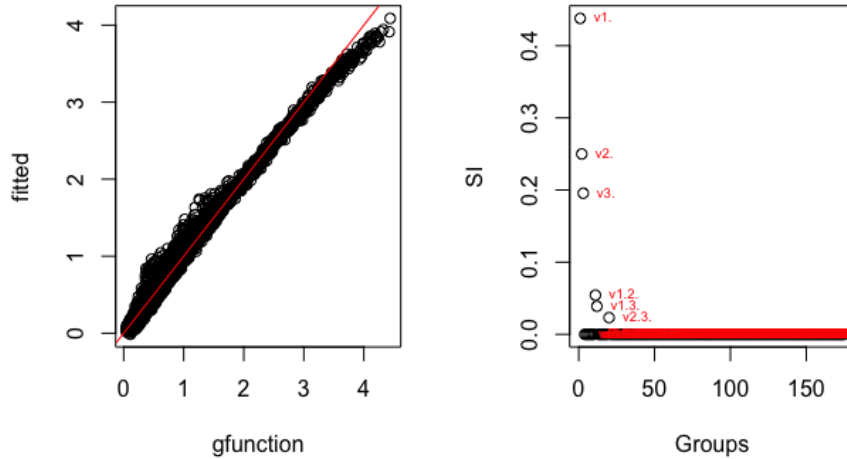


Figure 4.2: On the left, the RKHS meta-model versus the g-function is plotted. On the right, the empirical Sobol indices in the y axis and $v\text{Max}=175$ groups in the x axis are displayed.

4.6 Summary and discussion

An R package, called **RKHSMetaMod**, that estimates a meta-model of a complex model m , is proposed. This meta-model belongs to a reproducing kernel Hilbert space constructed as a direct sum of Hilbert spaces (Durrande et al. (2013)). The estimation of the meta-model is carried out via a penalized least squares minimization allowing both to select and estimate the terms in the Hoeffding decomposition, and therefore, to select the Sobol indices that are non-zero and estimate them (Huet and Taupin (2017)). This procedure makes it possible to estimate Sobol indices of high order, a point known to be difficult in practice.

Using the convex optimization tools, **RKHSMetaMod** package implements two optimization algorithms: the minimization of the RKHS ridge group sparse criterion (4.12) and the RKHS group lasso criterion (4.13). Both of these algorithms rely on the Gram matrices K_v , $v \in \mathcal{P}$ and their positive definiteness.

Currently, the package considers only uniformly distributed input variables. If one is interested by another distribution of the input variables, it suffices to modify the calculation of the kernels k_{0a} , $a = 1, \dots, d$ (see Equation (4.15)) in the function `calc_Kv` of this package (see Remark 4.3.1).

The available kernels in the **RKHSMetaMod** package are: linear kernel, quadratic kernel, brownian kernel, matern kernel and gaussian kernel (see Table 4.1). Regarding to the problem under study, one may consider another kernel and add it easily to the list of the kernels in the `calc_Kv` function. Indeed, the choice of different

kernels allows to consider different approximation spaces and choose the one that gives the best result.

For the large values of n and d the calculation and storage of the eigenvalues and the eigenvectors of all Gram matrices $K_v, v \in \mathcal{P}$ requires a lot of time and a very large amount of memory. In order to optimize the execution time and also the storage memory, except for a function that is written in R, all of the functions of **RKHSMetaMod** package are written using the efficient C++ libraries through **RcppEigen** and **RcppGSL** packages. These functions are then interfaced with the R environment in order to propose an user friendly package.

The performance of the package functions in terms of the predictive quality of the estimator and the estimation of the Sobol indices, is validated by a simulation study (see Examples 4.5.1-4.5.4).

The strategy of choosing the tuning parameters in this package is based on the minimization of the prediction error of the estimated meta-model, the prediction error being estimated using a testing dataset. The *best* estimator is selected in terms of the prediction quality, and the Sobol indices are deduced from it. If one is specially interested in the estimation of the Sobol indices, an alternative to our approach could be to calculate the tuning parameters which minimize the prediction error of the Sobol indices.

Appendix

4.A More technical details

Preliminary 4.A.1 For $F(x) = \|Ax\|$, where A is a symmetric matrix that not depends on x , the sub-differential of F at point x , denoted by $\partial F(x)$, is defined as follows:

$$\begin{aligned} \partial F(x) &= \left\{ \frac{A^2 x}{\|Ax\|} \right\} && \text{if } x \neq 0, \\ \partial F(x) &= \{w \in \mathbb{R}^n, \|A^{-1}w\| \leq 1\} && \text{if } x = 0. \end{aligned}$$

Preliminary 4.A.2 Let $F : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. we have the following first order optimality condition:

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} F(x) \Leftrightarrow 0 \in \partial F(\hat{x}).$$

This follows from the fact that $F(y) \geq F(\hat{x}) + \langle 0, y - \hat{x} \rangle$ for all $y \in \mathbb{R}^n$ in both cases (Giraud (2014)).

4.A.1 RKHS group lasso algorithm

We consider the minimization of the RKHS group lasso criterion given by,

$$C_g(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \mu_g \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\|.$$

We begin with the constant term f_0 . The ordinary first derivative of the function $C_g(f_0, \theta)$ at f_0 is equal to:

$$\frac{\partial C_g}{\partial f_0} = -2 \sum_{i=1}^n (Y_i - f_0) - \sum_{v \in \mathcal{P}} K_v \theta_v,$$

and therefore,

$$\hat{f}_0 = \frac{1}{n} \sum_{i=1}^n Y_i - \frac{1}{n} \sum_i \sum_v (K_v \theta_v)_i,$$

where $(K_v \theta_v)_i$ denotes the i -th component of $K_v \theta_v$.

Next step is to calculate,

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^n \times |\mathcal{P}|} C_g(f_0, \theta).$$

Since $C_g(f_0, \theta)$ is convex and separable, we use a block coordinate descent algorithm, group v by group v . In the following, we fix a group v , and we find the minimizer of $C_g(f_0, \theta)$ with respect to θ_v for given values of f_0 and θ_w , $w \neq v$. Set

$$C_{g,v}(f_0, \theta_v) = \|R_v - K_v \theta_v\|^2 + \sqrt{n} \mu_g \|K_v^{1/2} \theta_v\|,$$

where

$$R_v = Y - f_0 - \sum_{w \neq v} K_w \theta_w. \quad (4.19)$$

We aim to minimize $C_{g,v}(f_0, \theta_v)$ with respect to θ_v . Let $\partial C_{g,v}$ be the sub-differential of $C_{g,v}(f_0, \theta_v)$ with respect to θ_v :

$$\partial C_{g,v}(f_0, \theta) = \{-2K_v(R_v - K_v \theta_v) + \sqrt{n} \mu_g t_v : t_v \in \partial \|K_v^{1/2} \theta_v\|\}.$$

The first order optimality condition (see Preliminary (4.A.2)) ensures the existence of $\hat{t}_v \in \partial \|K_v^{1/2} \theta_v\|$ fulfilling,

$$-2K_v(R_v - K_v \theta_v) + \sqrt{n} \mu_g \hat{t}_v = 0. \quad (4.20)$$

Using the sub-differential definition (see Preliminary 4.A.1) we obtain,

$$\partial \|K_v^{1/2} \theta_v\| = \left\{ \frac{K_v \theta_v}{\|K_v^{1/2} \theta_v\|} \right\} \quad \text{if } \theta_v \neq 0,$$

and,

$$\partial \|K_v^{1/2} \theta_v\| = \{\hat{t}_v \in \mathbb{R}^n, \|K_v^{-1/2} \hat{t}_v\| \leq 1\} \quad \text{if } \theta_v = 0.$$

Let $\hat{\theta}_v$ be the minimizer of $C_{g,v}$. The sub-differential equations above give the two following cases:

Case 1. If $\hat{\theta}_v = 0$ then there exists $\hat{t}_v \in \mathbb{R}^n$ such that $\|K_v^{-1/2} \hat{t}_v\| \leq 1$ and it fulfils Equation (4.20):

$$2K_v R_v = \sqrt{n} \mu_g \hat{t}_v,$$

So, the necessary and sufficient condition for which the solution $\hat{\theta}_v = 0$ is the optimal one is:

$$\left\| \frac{2}{\sqrt{n}} K_v^{1/2} R_v \right\| \leq \mu_g.$$

Case 2. If $\hat{\theta}_v \neq 0$ then $\hat{t}_v = K_v \hat{\theta}_v / \|K_v^{1/2} \hat{\theta}_v\|$ and it fulfils Equation (4.20):

$$2K_v(R_v - K_v \hat{\theta}_v) = \sqrt{n} \mu_g \frac{K_v \hat{\theta}_v}{\|K_v^{1/2} \hat{\theta}_v\|}.$$

We obtain then,

$$\hat{\theta}_v = \left(K_v + \frac{\sqrt{n} \mu_g}{2 \|K_v^{1/2} \hat{\theta}_v\|} I_n \right)^{-1} R_v. \quad (4.21)$$

Since $\hat{\theta}_v$ appears in both sides of the Equation (4.21), a numerical procedure is needed:

Proposition 4.A.1 For $\rho > 0$ let $\theta(\rho) = (K_v + \rho I_n)^{-1} R_v$. There exists a non-zero solution to Equation (4.21) if and only if there exists $\rho > 0$ such that

$$\mu_g = \frac{2\rho}{\sqrt{n}} \|K_v^{1/2} \theta(\rho)\|. \quad (4.22)$$

Then $\hat{\theta}_v = \theta(\rho)$.

Proof If there exists a non-zero solution to Equation (4.21), then $\|K_v^{1/2}\widehat{\theta}_v\| \neq 0$ since K_v is positive definite. Take

$$\rho = \frac{\sqrt{n}\mu_g}{2\|K_v^{1/2}\widehat{\theta}_v\|},$$

then

$$\theta(\rho) = (K_v + \frac{\sqrt{n}\mu_g}{2\|K_v^{1/2}\widehat{\theta}_v\|}I_n)^{-1}R_v = \widehat{\theta}_v,$$

and, for such ρ Equation (4.22) is satisfied.

Conversely, if there exists $\rho > 0$ such that Equation (4.22) is satisfied, then $\|K_v^{1/2}\theta(\rho)\| \neq 0$ and,

$$\rho = \frac{\sqrt{n}\mu_g}{2\|K_v^{1/2}\theta(\rho)\|}.$$

Therefore,

$$\theta(\rho) = (K_v + \frac{\sqrt{n}\mu_g}{2\|K_v^{1/2}\theta(\rho)\|}I_n)^{-1}R_v,$$

which is Equation (4.21) calculated in $\widehat{\theta}_v = \theta(\rho)$. \square

Remark 4.A.1 Define $y(\rho) = 2\rho\|K_v^{1/2}\theta(\rho)\| - \sqrt{n}\mu_g$ with $\theta(\rho) = (K_v + \rho I_n)^{-1}R_v$, then $y(\rho) = 0$ has a unique solution, denoted $\widehat{\rho}$, which leads to calculate $\widehat{\theta}(\widehat{\rho})$.

Proof For $\rho = 0$ we have $y(0) = -\sqrt{n}\mu_g < 0$, since $\mu_g > 0$; and for $\rho \rightarrow +\infty$ we have $y(\rho) > 0$, since $\|K_v^{1/2}(\frac{K_v}{\rho} + I_n)^{-1}R_v\| \rightarrow \|K_v^{1/2}R_v\|$ and $\|2K_v^{1/2}R_v\| > \sqrt{n}\mu_g$.

Moreover, we have

$$\begin{aligned} y(\rho) &= 2\|(\frac{I_n}{\rho} + k_v^{-1})^{-1}k_v^{-1/2}R_v\| - \sqrt{n}\mu_g, \\ &= 2(X^T A^{-2}X)^{1/2} - \sqrt{n}\mu_g, \end{aligned}$$

where $A = (I_n/\rho + k_v^{-1})$ and $X = k_v^{-1/2}R_v$. The first derivative of $y(\rho)$ in ρ is obtained by,

$$\frac{\partial y(\rho)}{\partial \rho} = (X^T A^{-2}X)^{-1/2} \frac{\partial (X^T A^{-2}X)}{\partial \rho},$$

and,

$$\begin{aligned} \frac{\partial (X^T A^{-2}X)}{\partial \rho} &= X^T \frac{\partial (A^{-1})^2}{\partial \rho} X, \\ &= 2X^T A^{-1}(-A^{-1} \frac{\partial A}{\partial \rho} A^{-1})X, \\ &= \frac{2}{\rho^2} \|A^{-3/2}X\|. \end{aligned}$$

Finally, we get

$$\frac{\partial y(\rho)}{\partial \rho} = \frac{2\|(\frac{I_n}{\rho} + k_v^{-1})^{-3/2}k_v^{-1/2}R_v\|}{\rho^2\|(\frac{I_n}{\rho} + k_v^{-1})^{-1}k_v^{-1/2}R_v\|} > 0.$$

So $y(\rho)$ is an increasing function of ρ , and the proof is complete. \square

In order to calculate ρ and so $\hat{\theta}_v = \theta(\rho)$ we use Algorithm 4 which is a part of the RKHS group lasso Algorithm 1 when $\hat{\theta}_v \neq 0$.

Algorithm 4 Algorithm to find ρ as well as $\hat{\theta}_v$

```

1: if  $\hat{\theta}_{\text{old}} = 0$  then  $\triangleright \hat{\theta}_{\text{old}}$  is  $\hat{\theta}_v$  computed in the previous step of the RKHS group
   lasso algorithm.
2:   Set  $\rho \leftarrow 1$  and calculate  $y(\rho)$ 
3:   if  $y(\rho) > 0$  then
4:     Find  $\hat{\rho}$  that minimizes  $y(\rho)$  on the interval  $[0, 1]$ 
5:   else
6:     repeat
7:       Set  $\rho \leftarrow \rho \times 10$  and calculate  $y(\rho)$ 
8:     until  $y(\rho) > 0$ 
9:     Find  $\hat{\rho}$  that minimizes  $y(\rho)$  on the interval  $[\rho/10, \rho]$ 
10:  end if
11: else
12:   Set  $\rho \leftarrow \frac{\sqrt{n\mu_g}}{2\|K_v^{1/2}\hat{\theta}_{\text{old}}\|}$  and calculate  $y(\rho)$ 
13:   if  $y(\rho) > 0$  then
14:     repeat
15:       Set  $\rho \leftarrow \rho/10$  and calculate  $y(\rho)$ 
16:     until  $y(\rho) < 0$ 
17:     Find  $\hat{\rho}$  that minimizes  $y(\rho)$  on the interval  $[\rho, \rho \times 10]$ 
18:   else
19:     repeat
20:       Set  $\rho \leftarrow \rho \times 10$  and calculate  $y(\rho)$ 
21:     until  $y(\rho) > 0$ 
22:     Find  $\hat{\rho}$  that minimizes  $y(\rho)$  on the interval  $[\rho/10, \rho]$ 
23:   end if
24: end if
25: calculate  $\hat{\theta}_v = \theta(\hat{\rho})$ 

```

4.A.2 RKHS ridge group sparse algorithm

We consider the minimization of the RKHS ridge group sparse criterion:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \gamma \sum_{v \in \mathcal{P}} \|K_v \theta_v\| + n \mu \sum_{v \in \mathcal{P}} \|K_v^{1/2} \theta_v\|.$$

The constant term f_0 is estimated as in the RKHS group lasso algorithm. In order to calculate $\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^n \times |\mathcal{P}|} C(f_0, \theta)$, we use once again the block coordinate descent algorithm group v by group v . In the following, we fix a group v , and we find the minimizer of $C(f_0, \theta)$ with respect to θ_v for given values of f_0 and θ_w , $w \neq v$. We aim at minimizing with respect to θ_v ,

$$C_v(f_0, \theta_v) = \|R_v - K_v \theta_v\|^2 + \sqrt{n} \gamma \|K_v \theta_v\| + n \mu \|K_v^{1/2} \theta_v\|,$$

where R_v is defined by (4.19).

Let ∂C_v be the sub-differential of $C_v(f_0, \theta_v)$ with respect to θ_v ,

$$\partial C_v = \{-2K_v(R_v - K_v \theta_v) + \sqrt{n} \gamma s_v + n \mu t_v : s_v \in \partial \|K_v \theta_v\|, t_v \in \partial \|K_v^{1/2} \theta_v\|\},$$

According to the first order optimality condition (see Preliminary 4.A.2), we know that there exists $\hat{s}_v \in \partial \|K_v \theta_v\|$ and $\hat{t}_v \in \partial \|K_v^{1/2} \theta_v\|$ such that,

$$-2K_v(R_v - K_v \theta_v) + \sqrt{n} \gamma \hat{s}_v + n \mu \hat{t}_v = 0. \quad (4.23)$$

The sub-differential definition (see Preliminary 4.A.1) gives,

$$\{\partial \|K_v^{1/2} \theta_v\| = \left\{ \frac{K_v \theta_v}{\|K_v^{1/2} \theta_v\|} \right\}, \partial \|K_v \theta_v\| = \left\{ \frac{K_v^2 \theta_v}{\|K_v \theta_v\|} \right\}\} \quad \text{if } \theta_v \neq 0,$$

and,

$$\{\partial \|K_v^{1/2} \theta_v\| = \{\hat{t}_v \in \mathbb{R}^n, \|K_v^{-1/2} \hat{t}_v\| \leq 1\}, \partial \|K_v \theta_v\| = \{\hat{s}_v \in \mathbb{R}^n, \|K_v^{-1} \hat{s}_v\| \leq 1\}\} \quad \text{if } \theta_v = 0.$$

Let $\hat{\theta}_v$ be the minimizer of the $C_v(f_0, \theta_v)$. Using the sub-differential equations above, the estimator $\hat{\theta}_v$, $v \in \mathcal{P}$ is obtained following two cases below:

Case 1. If $\hat{\theta}_v = 0$ then there exists $\hat{s}_v \in \mathbb{R}^n$ such that $\|K_v^{-1} \hat{s}_v\| \leq 1$ and it fulfils Equation (4.23):

$$2K_v R_v - n \mu \hat{t}_v = \sqrt{n} \gamma \hat{s}_v,$$

with $\hat{t}_v \in \mathbb{R}^n$, $\|K_v^{-1/2} \hat{t}_v\| \leq 1$. Set

$$J(\hat{t}_v) = \|2R_v - n \mu K_v^{-1} \hat{t}_v\|,$$

and,

$$J^* = \arg \min_{\hat{t}_v \in \mathbb{R}^n} \{J(\hat{t}_v), \text{ such that } \|K_v^{-1/2} \hat{t}_v\| \leq 1\}.$$

Then the solution to Equation (4.23) is zero if and only if $J^* \leq \gamma$.

Case 2. If $\hat{\theta}_v \neq 0$ then we have $\hat{s}_v = K_v^2 \hat{\theta}_v / \|K_v \hat{\theta}_v\|$, and $\hat{t}_v = K_v \hat{\theta}_v / \|K_v^{1/2} \hat{\theta}_v\|$ fulfilling Equation (4.23):

$$2K_v(R_v - K_v \hat{\theta}_v) = \sqrt{n} \gamma \frac{K_v^2 \hat{\theta}_v}{\|K_v \hat{\theta}_v\|_2} + n \mu \frac{K_v \hat{\theta}_v}{\|K_v^{1/2} \hat{\theta}_v\|},$$

that is,

$$\hat{\theta}_v = \left(K_v + \frac{\sqrt{n} \gamma}{2 \|K_v \hat{\theta}_v\|} K_v + \frac{n \mu}{2 \|K_v^{1/2} \hat{\theta}_v\|} I_n \right)^{-1} R_v \quad \text{if } \hat{\theta}_v \neq 0.$$

In this case the calculation of $\hat{\theta}_v$ needs a numerical algorithm which is explained in [Huet and Taupin \(2017\)](#).

Package ‘RKHSMetaMod’

Type Package

Title Ridge Group Sparse Optimization Problem for Estimation of Meta-Models Based on Reproducing Kernel Hilbert Space

Version 1.0

Date 2019-06-17

Author Halaleh Kamari

Maintainer Halaleh Kamari <halaleh.kamari@univ-evry.fr>

Description It estimates the Hoeffding decomposition of a complex function by solving ridge group sparse optimization problem based on a Reproducing Kernel Hilbert Space, and approximates its Sobol Indices.

License GPL (>=2.0)

Imports Rcpp (>= 1.0.0)

LinkingTo Rcpp, RcppEigen, RcppGSL

Contents

A.1	calc_Kv function	157
A.2	grplasso_q function	159
A.3	mu_max function	161
A.4	pen_MetMod function	163
A.5	PredErr function	166
A.6	RKHSgrplasso function	167
A.7	RKHSMetMod function	169
A.8	RKHSMetMod_qmax function	172
A.9	SI_emp function	175

Index

RKHSMetaMod-package	<i>Produces a sequence of meta-models that are the solutions of the RKHS ridge group sparse or the RKHS group lasso optimization problems.</i>
---------------------	--

Description

Estimates a meta-model that approximates the Hoeffding decomposition of a complex model m by solving the ridge group sparse (or group lasso) optimization prob-

lem based on a Reproducing Kernel Hilbert Space (RKHS). The model m depends on d input variables $X = (X_1, \dots, X_d)$ that are independent and uniformly distributed on $[0, 1]^d$. This model m from \mathbb{R}^d to \mathbb{R} may be a known model that can be calculated in all points of X , or it may be an unknown regression model defined as follows:

$$Y = m(X) + \sigma\varepsilon, \quad \sigma > 0,$$

where the error ε is assumed to be centered with a finite variance, i.e. $E(\varepsilon) = 0$ and $\text{var}(\varepsilon) < \infty$.

Let \mathcal{P} be the set of parts of $\{1, \dots, d\}$ with dimension 1 to d . The RKHS ridge group sparse criterion is defined by:

$$C(f_0, \theta) = \|Y - f_0 I_n - \sum_{v \in \mathcal{P}} K_v \theta_v\|^2 + \sqrt{n} \gamma \sum_{v \in \mathcal{P}} \gamma'_v \|K_v \theta_v\| + n \mu \sum_{v \in \mathcal{P}} \mu'_v \|K_v^{1/2} \theta_v\|,$$

where for all $v \in \mathcal{P}$, γ'_v and μ'_v are the vector of weights that should be chosen suitably, and K_v are the Gram matrices associated with a chosen reproducing kernel.

The RKHS group lasso criterion is obtained by setting $\gamma = 0$ in the RKHS ridge group sparse criterion above. The RKHS group lasso penalty parameter is denoted by $\mu_g = \sqrt{n} \mu$.

For each pair of the penalty parameters (μ, γ) in the RKHS ridge group sparse criterion, one estimator, called RKHS meta-model, is calculated. For a given value $D_{\max} \in \mathbb{N}$, the RKHS meta-model \hat{f} has an additive representation including the variables and interactions between them of order maximum equal to D_{\max} :

$$\hat{f} = f_0 + \sum_{v \in \mathcal{P}_{D_{\max}}} f_v,$$

where f_0 is a constant, and $\mathcal{P}_{D_{\max}}$ is the set of parts of $\{1, \dots, d\}$ with dimensions 1 to D_{\max} and cardinality equal to v_{\max} :

$$v_{\max} = \sum_{j=1}^{D_{\max}} \binom{d}{j}.$$

For a given grid of values of the tuning parameters (μ, γ) a sequence of the RKHS meta-models are produced by minimizing the RKHS ridge group sparse criterion (if $\gamma \neq 0$) or the RKHS group lasso criterion (if $\gamma = 0$). These meta-models are evaluated using a testing dataset. That is, the prediction error is calculated for each RKHS meta-model, and the one with the minimum prediction error is the *best* estimator for the true model m . This package provides a function that estimates the empirical Sobol indices of the obtained RKHS meta-models. The estimators of the Sobol indices of m are deduced from the *best* RKHS meta-model.

Details

Package: RKHSMetaMod
 Type: Package
 Version: 1.0
 Date: 2019-06-17
 License: GPL (>=2.0)

Author(s)

Halaleh Kamari

Maintainer: <halaleh.kamari@univ-evry.fr>

References

Kamari, H., Huet, S., Taupin, M.-L. (2019) RKHSMetaMod : An R package to estimate the Hoeffding decomposition of a complex model by solving RKHS ridge group sparse optimization problem. arXiv:1905.13695 [stat.ML].

See Also[RKHSMetaMod](#)**Examples**

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
frc <- c(10,100)
gamma <- c(.5,.01,.001,0)
result <- RKHSMetaMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
```

A.1 `calc_Kv` function

<code>calc_Kv</code>	<i>Function to calculate the eigenvalues and eigenvectors of the Gram matrices K_v, $v \in \mathcal{P}_{D_{\max}}$.</i>
----------------------	---

Description

For a given value of D_{\max} this function calculates the Gram matrices K_v for $v \in \mathcal{P}_{D_{\max}}$, and returns their associated eigenvalues and eigenvectors. The calculated Gram matrices may be not positive definite. The option "correction" of this function allows to replace the matrices K_v that are not positive definite by their

"nearest positive definite" matrices.

Usage

`calc_Kv(X, kernel, Dmax, correction, verbose, tol)`

Arguments

<code>X</code>	Matrix of observations with n rows and d columns.
<code>kernel</code>	Character, the type of the reproducing kernel: matern (matern kernel), brownian (brownian kernel), gaussian (gaussian kernel), linear (linear kernel), quad (quadratic kernel).
<code>Dmax</code>	Integer, between 1 and d , indicates the order of interactions considered in the meta-model: Dmax= 1 is used to consider only the main effects, Dmax= 2 to include the main effects and the interactions of order 2, ...
<code>correction</code>	Logical, if TRUE, the program makes the correction to the matrices K_v that are not positive definite (see details). Set as TRUE by default.
<code>verbose</code>	Logical, if TRUE, the group v for which the correction is done is printed. Set as TRUE by default.
<code>tol</code>	Scalar, used if correction is TRUE. For each matrix K_v if $\lambda_{min} < \lambda_{max} \times \text{tol}$, then the correction to K_v is done (see details). Set as $1e^{-8}$ by default.

Details

Let $\lambda_{v,i}, i = 1, \dots, n$ be the eigenvalues associated with matrix K_v . Set $\lambda_{max} = \max_i \lambda_{v,i}$ and $\lambda_{min} = \min_i \lambda_{v,i}$. The eigenvalues of K_v that is not positive definite are replaced by $\lambda_{v,i} + \text{epsilon}$, with $\text{epsilon} = \lambda_{max} \times \text{tol}$. The value of tol depends on the type of the kernel and it is chosen small.

Value

List of two components "names.Grp" and "kv":

<code>names.Grp</code>	Vector of size vMax, indicates the name of groups included in the meta-model.
<code>kv</code>	List of vMax components with the same names as the vector names.Grp. Each element of the list is a list of two components "Evalues" and "Q":
<code>Evalues</code>	Vector of size n , eigenvalues of each Gram matrix K_v .
<code>Q</code>	Matrix with n rows and n columns, eigenvectors of each Gram matrix K_v .

Note: Note.

Author(s)

Halaleh Kamari

References

References.

See Also

RKHSMetaMod

Examples

```

d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
Kv <- calc_Kv(X, kernel, Dmax)
names <- Kv$names.Grp
Eigen.val1 <- Kv$kv$V1.$Evalues
Eigen.vec1 <- Kv$kv$V1.$Q

```

A.2 `grplasso_q` function

`grplasso_q` *Function to fit a solution with q active groups of the RKHS group lasso optimization problem.*

Description

This function determines the value $\mu_q(q)$, for which the number of active groups in the solution of the RKHS group lasso problem is equal to q , and returns the RKHS meta-model associated with $\mu_q(q)$.

Usage

```
grplasso_q(Y, Kv, q, rat, Num)
```

Arguments

Y	Vector of response observations of size n .
Kv	List of eigenvalues and eigenvectors of positive definite Gram matrices K_v and their associated group names. It should have the same format as the output of the function <code>calc_Kv</code> (see details).
q	Integer, the number of active groups in the obtained solution.
rat	Positive scalar, used to restrict the minimum value of μ_g , to be evaluated in the RKHS group lasso algorithm, $\mu_{min} = \mu_{max}/rat$. The value μ_{max} is calculated inside the program, see function <code>mu_max</code> .
Num	Integer, used to restrict the number of different values of the penalty parameter μ_g to be evaluated in the RKHS group lasso algorithm, until it achieves $\mu_g(q)$: for Num= 1 the program is done for 3 values of μ_g , $\mu_1 = (\mu_{min} + \mu_{max})/2$, $\mu_2 = (\mu_{min} + \mu_1)/2$ or $\mu_2 = (\mu_1 + \mu_{max})/2$ depending on the value of q associated with μ_1 , $\mu_3 = \mu_{min}$.

Details

Input Kv should contain the eigenvalues and eigenvectors of positive definite Gram matrices K_v . It is necessary to set input "correction" in the function `calc_Kv` equal to "TRUE".

Value

List of 4 components: "mus", "qs", "mu", "res":

mus	Vector, values of the evaluated penalty parameters μ_g in the RKHS group lasso algorithm until it achieves $\mu_g(q)$.
qs	Vector, number of active groups associated with each value of μ_g in mus.
mu	Scalar, value of $\mu_g(q)$.
res	An RKHS group lasso object:
intercept	Scalar, estimated value of intercept.
teta	Matrix with vMax rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, vMax$.
fit.v	Matrix with n rows and vMax columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
fitted	Vector of size n , indicates the estimator of m .
Norm.H	Vector of size vMax, estimated values of the penalty norm.
supp	Vector of active groups.
Nsupp	Vector of the names of the active groups.
SCR	Scalar, equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
crit	Scalar, indicates the value of the penalized criterion.
MaxIter	Integer, number of iterations until convergence is reached.
convergence	TRUE or FALSE. Indicates whether the algorithm has converged or not.
RelDiffCrit	Scalar, value of the first convergence criterion at the last iteration, $\frac{crit_{lastIter} - crit_{lastIter-1}}{crit_{lastIter-1}}$.
RelDiffPar	Scalar, value of the second convergence criterion at the last iteration, $\ \frac{\theta_{lastIter} - \theta_{lastIter-1}}{\theta_{lastIter-1}}\ ^2$.

Note:

Note.

Author(s)

Halaleh Kamari

References

References.

See Also

[calc_Kv](#), [mu_max](#)

Examples

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
Kv <- calc_Kv(X, kernel, Dmax, TRUE, TRUE)
result <- grplasso_q(Y,Kv,5,100 ,Num=10)
result$mu
result$res$Nsupp
```

A.3 `mu_max` function

<code>mu_max</code>	<i>Function to find the maximal value of the penalty parameter in the RKHS group lasso optimization problem.</i>
---------------------	--

Description

Calculates the value of the penalty parameter in the RKHS group lasso optimization problem when the first penalized parameter group enters the model.

Usage

```
mu_max(Y, matZ)
```

Arguments

Y Vector of response observations of size n .
 List of `vMax` components. Each component includes the eigenvalues and eigenvectors of the positive definite Gram matrices

matZ $K_v, v = 1, \dots, vMax$. It should have the same format as the output "kv" of the function `calc_Kv`.

Details

For more details about the maximal value of the penalty parameter in the ordinary group lasso algorithm see Meier et al. (2008).

Value

An object of type numeric is returned.

Note:

Note.

Author(s)

Halaleh Kamari

References

Meier, L. Van de Geer, S. and Bühlmann, P. (2008) The group lasso for logistic regression, Eidgenössische Technische Hochschule, Zurich, Switzerland. J. R. Statist. Soc. B (2008) 70, Part 1, pp. 53-71.

See Also

`calc_Kv`

Examples

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
Kv <- calc_Kv(X, kernel, Dmax, TRUE,TRUE)
matZ <- Kv$kv
mumax <- mu_max(Y, matZ)
mumax
```

A.4 `pen_MetMod` function

<code>pen_MetMod</code>	<i>Function to fit a solution of the RKHS ridge group sparse optimization problem.</i>
-------------------------	--

Description

This function produces a sequence of the RKHS meta-models associated with a given grid of values of the tuning parameters μ, γ . Each RKHS meta-model in the sequence is the solution to the RKHS ridge group sparse optimization problem associated with a pair of values of (μ, γ) in the grid of values of μ, γ .

Usage

```
pen_MetMod(Y, Kv, gamma, mu, resg, gama_v, mu_v, maxIter, verbose, calcStwo)
```

Arguments

<code>Y</code>	Vector of response observations of size n .
<code>Kv</code>	List, includes the eigenvalues and eigenvectors of the positive definite Gram matrices $K_v, v = 1, \dots, vMax$ and their associated group names. It should have the same format as the output of the function <code>calc_Kv</code> (see details).
<code>gamma</code>	Vector of positive scalars. Values of the penalty parameter γ in decreasing order.
<code>mu</code>	Vector of positive scalars. Values of the penalty parameter μ in decreasing order.
<code>resg</code>	List of initial parameters, includes the <code>RKHSgrplasso</code> objects for each value of the penalty parameter μ .
<code>gama_v</code>	Scalar zero or vector of <code>vMax</code> positive scalars, considered as weights for the ridge penalty. Set to zero, to consider no weights, i.e. all weights equal to 1.
<code>mu_v</code>	Scalar zero or a vector with <code>vMax</code> scalars, considered as weights of sparse group penalty. Set to zero, to consider no weights, i.e. all weights equal to 1.
<code>maxIter</code>	Integer, shows the maximum number of loops through initial active groups at the first step and maximum number of loops through all groups at the second step. Set as 1000 by default.
<code>verbose</code>	Logical, if TRUE, for each pair of penalty parameters (μ, γ) it prints: the number of current iteration, active groups and convergence criteria. Set as FALSE by default.
<code>calcStwo</code>	Logical, if TRUE, the program does a second step after convergence: the algorithm is done over all groups by taking the estimated parameters at the first step as initial values. Set as FALSE by default.

Details

Input `Kv` should contain the eigenvalues and eigenvectors of positive definite Gram

matrices K_v . It is necessary to set input "correction" in the function `calc_Kv` equal to "TRUE".

Value

List of l components, with l equals to the number of pairs of the penalty parameters (μ, γ) . Each component of the list is a list of 3 components "mu", "gamma" and "Meta-Model":

<code>mu</code>	Positive scalar, an element of the input vector mu associated with the estimated Meta-Model.
<code>gamma</code>	Positive scalar, an element of the input vector gamma associated with the estimated Meta-Model.
<code>Meta-Model</code>	Estimated meta-model associated with penalty parameters mu and gamma. List of 16 components:
<code>intercept</code>	Scalar, estimated value of intercept.
<code>teta</code>	Matrix with vMax rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, vMax$.
<code>fit.v</code>	Matrix with n rows and vMax columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
<code>fitted</code>	Vector of size n, indicates the estimator of m.
<code>Norm.n</code>	Vector of size vMax, estimated values for the ridge penalty norm.
<code>Norm.H</code>	Vector of size vMax, estimated values for the group sparse penalty norm.
<code>supp</code>	Vector of active groups.
<code>Nsupp</code>	Vector of the names of the active groups.
<code>SCR</code>	Scalar equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
<code>crit</code>	Scalar indicates the value of the penalized criterion.
<code>gamma.v</code>	Vector of size vMax, coefficients of the ridge penalty norm, $\sqrt{n} \gamma \times \text{gama}_v$.
<code>mu.v</code>	Vector of size vMax, coefficients of the group sparse penalty norm, $n \mu \times \text{mu}_v$.
<code>iter</code>	List of three components if calcStwo=TRUE (two components if calcStwo=FALSE): maxIter, number of iterations until convergence is reached at first step and the number of iterations until convergence is reached at second step (maxIter, and the number of iterations until convergence is reached at first step).
<code>convergence</code>	TRUE or FALSE. Indicates whether the algorithm has converged or not.
<code>RelDiffCrit</code>	List of two components if calcStwo=TRUE (one component if calcStwo=FALSE): value of convergence criterion at the last iteration of each step, $\left\ \frac{\theta_{lastIter} - \theta_{lastIter-1}}{\theta_{lastIter-1}} \right\ ^2$.
<code>RelDiffPar</code>	List of two components if calcStwo=TRUE (one component if calcStwo=FALSE): value of convergence criterion at the last iteration, $\frac{crit_{lastIter} - crit_{lastIter-1}}{crit_{lastIter-1}}$ of each step.

Note:

For more details about the algorithm see Huet and Taupin (2017).

Author(s)

Halaleh Kamari

References

Huet, S. and Taupin, M. L. (2017). Metamodel construction for sensitivity analysis. ESAIM: Procs 60, 27-69.

See Also

[calc_Kv](#), [RKHSgrplasso](#)

Examples

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
Kv <- calc_Kv(X, kernel, Dmax, TRUE,TRUE, tol = 1e-8)
vMax <- length(Kv$names.Grp)
matZ <- Kv$kv
mumax <- mu_max(Y, matZ)
mug1 <- mumax/10
mug2 <- mumax/100
gr1 <- RKHSgrplasso(Y,Kv, mug1)
gr2 <- RKHSgrplasso(Y,Kv, mug2)
gamma <- c(.5,.01,.001)
rescaling the penalty parameter
mu <- c(mug1/sqrt(n),mug2/sqrt(n))
resg<-list(gr1,gr2)
res <- pen_MetMod(Y,Kv,gamma,mu,resg,0,0)
l <- length(res)
for(i in 1:l)print(res[[i]]$mu)
for(i in 1:l)print(res[[i]]$gamma)
for(i in 1:l)print(res[[i]]$'Meta-Model'$Nsupp)
gama_v <- rep(1,vMax)
mu_v <- rep(1,vMax)
res.w <- pen_MetMod(Y,Kv,gamma,mu,resg,gama_v,mu_v)
for(i in 1:l)print(res.w[[i]]$'Meta-Model'$Nsupp)
```


A.5 PredErr function

`PredErr` *Function to calculate the prediction error.*

Description

Computes the prediction error by considering a testing dataset.

Usage

`PredErr(X, XT, YT, mu, gamma, res, kernel, Dmax)`

Arguments

<code>X</code>	Matrix of observations with n rows and d columns.
<code>XT</code>	Matrix of observations of the testing dataset with n^{test} rows and d columns.
<code>YT</code>	Vector of response observations of testing dataset of size n^{test} .
<code>mu</code>	Vector of positive scalars. Values of the group sparse penalty parameter in decreasing order. See function RKHSMetMod .
<code>gamma</code>	Vector of positive scalars. Values of the ridge penalty parameter in decreasing order. See function RKHSMetMod .
<code>res</code>	List, includes a sequence of estimated meta-models for the learning dataset, using RKHS ridge group sparse or RKHS group lasso algorithm, associated with the penalty parameters <code>mu</code> and <code>gamma</code> . It should have the same format as the output of one of the functions: pen_MetMod , RKHSMetMod or RKHSMetMod_qmax .
<code>kernel</code>	Character, shows the type of the reproducing kernel: <code>matern</code> , <code>brownian</code> , <code>gaussian</code> , <code>linear</code> , <code>quad</code> . The same kernel should be chosen as the one used for the learning dataset. See function calc_Kv .
<code>Dmax</code>	Integer between 1 and d . The same <code>Dmax</code> should be chosen as the one used for learning dataset. See function calc_Kv .

Details

Details.

Value

Matrix of the prediction errors is returned. Each element of the matrix is the obtained prediction error associated with one RKHS meta-model in "res".

Note:

Note.

Author(s)

Halaleh Kamari

References

References.

See Also

[calc_Kv](#), [pen_MetMod](#), [RKHSMetMod](#), [RKHSMetMod_qmax](#)

Examples

```
d <- 3
n <- 50
nT <- 50
library(lhs)
X <- maximinLHS(n, d)
XT <- maximinLHS(nT, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
FT <- 1;for (a in 1:d) FT <- FT*(abs(4*XT[,a]-2)+c[a])/(1+c[a])
sigma <- 0.2
epsilon <- rnorm(n,0,1);Y <- F + sigma*epsilon
epsilonT <- rnorm(nT,0,1);YT <- FT + sigma*epsilonT
Dmax <- 3
kernel <- "matern"
frc <- c(10,100)
gamma=c(.5,.01,.001)
res <- RKHSMetMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
mu <- vector()
l <- length(gamma)
for(i in 1:length(frc))mu[i]=res[[i-1]*l+1]$mu
error <- PredErr(X,XT, YT,mu,gamma, res, kernel,Dmax)
error
```

A.6 RKHSgrplasso function

RKHSgrplasso	<i>Function to fit a solution of an RKHS group lasso optimization problem.</i>
--------------	--

Description

For a given value of the tuning parameter μ_g , this function fits the solution to the RKHS group lasso optimization problem.

Usage

```
RKHSgrplasso(Y, Kv, mu, maxIter, verbose)
```

Arguments

Y	Vector of response observations of size n .
Kv	List, includes the eigenvalues and eigenvectors of the positive definite Gram matrices $K_v, v = 1, \dots, vMax$ and their associated group names. It should have the same format as the output of the function <code>calc_Kv</code> (see details).
mu	Positive scalar, value of the penalty parameter μ_g in the RKHS group lasso problem.
maxIter	Integer, shows the maximum number of loops through all groups. Set as 1000 by default.
verbose	Logical, if TRUE, prints: the number of current iteration, active groups and convergence criteria. Set as FALSE by default.

Details

Input Kv should contain the eigenvalues and eigenvectors of positive definite Gram matrices K_v . It is necessary to set input "correction" in the function `calc_Kv` equal to "TRUE".

For more details about the ordinary group lasso algorithm see Meier et al. (2008).

Value

Estimated RKHS meta-model, list with 13 components:

intercept	Scalar, estimated value of intercept.
teta	Matrix with $vMax$ rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, vMax$.
fit.v	Matrix with n rows and $vMax$ columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
fitted	Vector of size n , indicates the estimator of m .
Norm.H	Vector of size $vMax$, estimated values of the penalty norm.
supp	Vector of active groups.
Nsupp	Vector of the names of the active groups.
SCR	Scalar equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
crit	Scalar indicates the value of the penalized criterion.
MaxIter	Integer, number of iterations until convergence is reached.
convergence	TRUE or FALSE. Indicates whether the algorithm has converged or not.
RelDiffCrit	Scalar, value of the first convergence criterion at the last iteration, $\frac{crit_{lastIter} - crit_{lastIter-1}}{crit_{lastIter-1}}$.
RelDiffPar	Scalar, value of the second convergence criterion at the last iteration, $\ \frac{\theta_{lastIter} - \theta_{lastIter-1}}{\theta_{lastIter-1}}\ ^2$.

Note:

Note.

Author(s)

Halaleh Kamari

References

Meier, L. Van de Geer, S. and Bühlmann, P. (2008) The group lasso for logistic regression, Eidgenössische Technische Hochschule, Zurich, Switzerland. J. R. Statist. Soc. B (2008) 70, Part 1, pp. 53-71.

See Also

[calc_Kv](#)

Examples

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
Kv <- calc_Kv(X, kernel, Dmax, TRUE, TRUE)
matZ <- Kv$kv
mumax <- mu_max(Y, matZ)
mug <- mumax/10
gr <- RKHSgrplasso(Y,Kv, mug , 1000, FALSE)
gr$Nsupp
```

A.7 RKHSMetMod function

RKHSMetMod

Function to produce a sequence of the RKHS meta-models that are the solutions of the RKHS ridge group sparse or the RKHS group lasso optimization problems.

Description

For a given value of D_{\max} and a chosen reproducing kernel, this function calculates the Gram matrices K_v , $v \in \mathcal{P}_{D_{\max}}$, and produces a sequence of estimators \hat{f} associated with a given grid of values of tuning parameters μ, γ , i.e. the solutions to the RKHS ridge group sparse (if $\gamma \neq 0$) or the RKHS group lasso problem (if $\gamma = 0$).

Usage

```
RKHSMetMod(Y, X, kernel, Dmax, gamma, frc, verbose)
```

Arguments

Y	Vector of response observations of size n .
X	Matrix of observations with n rows and d columns.
kernel	Character, the type of the reproducing kernel: matern (matern kernel), brownian (brownian kernel), gaussian (gaussian kernel), linear (linear kernel), quad (quadratic kernel).
Dmax	Integer, between 1 and d , indicates the order of interactions considered in the meta-model: Dmax= 1 is used to consider only the main effects, Dmax= 2 to include the main effects and the interactions of order 2,
gamma	Vector of non negative scalars, values of the penalty parameter γ in decreasing order. If $\gamma = 0$ the function solves an RKHS group lasso problem and for $\gamma > 0$ it solves an RKHS ridge group sparse problem.
frc	Vector of positive scalars. Each element of the vector sets a value to the penalty parameter μ , $\mu = \mu_{max}/(\sqrt{n} \times frc)$. The value μ_{max} is calculated by the program. See the function mu_max .
verbose	Logical, if TRUE, prints: the group v for which the correction of Gram matrix K_v is done, and for each pair of the penalty parameters (μ, γ) : the number of current iteration, active groups and convergence criteria. Set as FALSE by default.

Details

Details.

Value

List of l components, with l equals to the number of pairs of the penalty parameters (μ, γ) . Each component of the list is a list of 3 components "mu", "gamma" and "Meta-Model":

<code>mu</code>	Positive scalar, penalty parameter μ associated with the estimated Meta-Model.
<code>gamma</code>	Positive scalar, an element of the input vector gamma associated with the estimated Meta-Model.
<code>Meta-Model</code>	An RKHS ridge group sparse or RKHS group lasso object associated with the penalty parameters mu and gamma:
<code>intercept</code>	Scalar, estimated value of intercept.
<code>teta</code>	Matrix with <code>vMax</code> rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, \text{vMax}$.
<code>fit.v</code>	Matrix with n rows and <code>vMax</code> columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
<code>fitted</code>	Vector of size n , indicates the estimator of m .
<code>Norm.n</code>	Vector of size <code>vMax</code> , estimated values for the ridge penalty norm.
<code>Norm.H</code>	Vector of size <code>vMax</code> , estimated values for the group sparse penalty norm.
<code>supp</code>	Vector of active groups.
<code>Nsupp</code>	Vector of the names of the active groups.
<code>SCR</code>	Scalar equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
<code>crit</code>	Scalar indicates the value of the penalized criterion.
<code>gamma.v</code>	Vector of size <code>vMax</code> , coefficients of the ridge penalty norm, $\sqrt{n}\gamma \times \text{gama}_v$.
<code>mu.v</code>	Vector of size <code>vMax</code> , coefficients of the group sparse penalty norm, $n\mu \times \text{mu}_v$.
<code>iter</code>	List of two components: <code>maxIter</code> , and the number of iterations until the convergence is achieved.
<code>convergence</code>	TRUE or FALSE. Indicates whether the algorithm has converged or not.
<code>RelDiffCrit</code>	Scalar, value of the first convergence criterion at the last iteration, $\left\ \frac{\theta_{\text{lastIter}} - \theta_{\text{lastIter}-1}}{\theta_{\text{lastIter}-1}} \right\ ^2$.
<code>RelDiffPar</code>	Scalar, value of the second convergence criterion at the last iteration, $\frac{\text{crit}_{\text{lastIter}} - \text{crit}_{\text{lastIter}-1}}{\text{crit}_{\text{lastIter}-1}}$.

Note:

Note

Author(s)

Halaleh Kamari

References

References.

See Also

[mu_max](#), [RKHSgrplasso](#), [pen_MetMod](#)

Examples

```

d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
frc <- c(10,100)
gamma <- c(.5,.01,.001,0)
result <- RKHSMetMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
l <- length(result)
for(i in 1:l)print(result[[i]]$mu)
for(i in 1:l)print(result[[i]]$gamma)
for(i in 1:l)print(result[[i]]$'Meta-Model'$Nsupp)

```

A.8 RKHSMetMod_qmax function

RKHSMetMod_qmax

Function to produce a sequence of the RKHS meta-models, with at most qmax active groups in their support. These meta-models are the solutions of the RKHS ridge group sparse or the RKHS group lasso optimization problems.

Description

For a given value of Dmax and a chosen reproducing kernel, this function calculates the Gram matrices K_v , $v \in \mathcal{P}_{Dmax}$, determines μ , denoted μ_{qmax} , for which the number of active groups in the RKHS group lasso solution is equal to $qmax$, and produces a sequence of the RKHS meta-models associated with the tuning parameter μ_{qmax} and a grid of values of the tuning parameter γ . All the RKHS meta-models produced by this function have at most $qmax$ active groups in their support.

Usage

```
RKHSMetMod_qmax(Y, X, kernel, Dmax, gamma, qmax, rat, Num, verbose)
```

Arguments

Y	Vector of response observations of size n .
X	Matrix of observations with n rows and d columns.
kernel	Character, indicates the type of the reproducing kernel: matern (matern kernel), brownian (brownian kernel), gaussian (gaussian kernel), linear (linear kernel), quad (quadratic kernel).
Dmax	Integer, between 1 and d , indicates the order of interactions considered in the meta-model: Dmax= 1 is used to consider only the main effects, Dmax= 2 to include the main effects and the interactions of order 2, . . .
gamma	Vector of non negative scalars, values of the penalty parameter γ in decreasing order. If $\gamma = 0$ the function solves an RKHS group lasso problem and for $\gamma > 0$ it solves an RKHS ridge group sparse problem.
qmax	Integer, shows the maximum number of active groups in the obtained solution.
rat	Positive scalar, to restrict the minimum value of μ considered in the algorithm, $\mu_{min} = \mu_{max}/(\sqrt{n} \times rat)$. The value μ_{max} is calculated inside the program, see function <code>mu_max</code> .
Num	Integer, it is used to restrict the number of different values of the penalty parameter μ to be evaluated in the RKHS group lasso algorithm until it achieves $\mu(qmax)$: for Num= 1 the program is done for 3 different values of μ , $\mu_1 = (\mu_{min} + \mu_{max})/2$, $\mu_2 = (\mu_{min} + \mu_1)/2$ or $\mu_2 = (\mu_1 + \mu_{max})/2$ depending on the number of active groups in the meta-model associated with μ_1 , $\mu_3 = \mu_{min}$.
verbose	Logical, if TRUE, prints: the group v for which the correction of Gram matrix K_v is done, and for each pair of (μ, γ) : the number of current iteration, active groups and convergence criteria. Set as FALSE by default.

Details

Details.

Value

List of three components "mus", "qs", and "MetaModel":

<code>mus</code>	Vector, values of the evaluated penalty parameters μ in the RKHS group lasso algorithm until it achieves $\mu(qmax)$.
<code>qs</code>	Vector, number of active groups associated with each element in <code>mus</code> .
<code>MetaModel</code>	List with the same length as the vector <code>gamma</code> . Each component of the list is a list of 3 components "mu", "gamma" and "Meta-Model":
<code>mu</code>	Scalar, the value $\mu(qmax)$.
<code>gamma</code>	Positive scalar, element of the input vector <code>gamma</code> associated with the estimated Meta-Model.
<code>Meta-Model</code>	An RKHS ridge group sparse or RKHS group lasso object associated with the penalty parameters <code>mu</code> and <code>gamma</code> :
<code>intercept</code>	Scalar, estimated value of intercept.
<code>teta</code>	Matrix with <code>vMax</code> rows and n columns. Each row of the matrix is the estimated vector θ_v for $v = 1, \dots, vMax$.
<code>fit.v</code>	Matrix with n rows and <code>vMax</code> columns. Each row of the matrix is the estimated value of $f_v = K_v \theta_v$.
<code>fitted</code>	Vector of size n , indicates the estimator of m .
<code>Norm.n</code>	Vector of size <code>vMax</code> , estimated values for the ridge penalty norm.
<code>Norm.H</code>	Vector of size <code>vMax</code> , estimated values for the group sparse penalty norm.
<code>supp</code>	Vector of active groups.
<code>Nsupp</code>	Vector of the names of the active groups.
<code>SCR</code>	Scalar equals to $\ Y - f_0 - \sum_v K_v \theta_v\ ^2$.
<code>crit</code>	Scalar indicates the value of the penalized criterion.
<code>gamma.v</code>	Vector of size <code>vMax</code> , coefficients of the ridge penalty norm, $\sqrt{n}\gamma \times \text{gamma}_v$.
<code>mu.v</code>	Vector of size <code>vMax</code> , coefficients of the group sparse penalty norm, $n\mu \times \text{mu}_v$.
<code>iter</code>	List of two components: <code>maxIter</code> , and the number of iterations until the convergence is achieved.
<code>convergence</code>	TRUE or FALSE. Indicates whether the algorithm has converged or not.
<code>RelDiffCrit</code>	Scalar, value of the first convergence criterion at the last iteration, $\left\ \frac{\theta_{lastIter} - \theta_{lastIter-1}}{\theta_{lastIter-1}} \right\ ^2$.
<code>RelDiffPar</code>	Scalar, value of the second convergence criterion at the last iteration, $\frac{crit_{lastIter} - crit_{lastIter-1}}{crit_{lastIter-1}}$.

Note:

For the case $\gamma = 0$ the outputs "mu" = μ_g and "Meta-Model" is the same as the one returned by the function [RKHSgrlasso](#).

Author(s)

Halaleh Kamari

References

Reference.

See Also

[mu_max](#), [RKHSgrplasso](#), [pen_MetMod](#), [grplasso_q](#)

Examples

```
d <- 3
n <- 50
library(lhs)
X <- maximinLHS(n, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
epsilon <- rnorm(n,0,1);sigma <- 0.2
Y <- F + sigma*epsilon
Dmax <- 3
kernel <- "matern"
gamma <- c(.5,.01,.001,0)
Num <- 10
rat <- 100
qmax <- 4
result <- RKHSMetMod_qmax(Y, X, kernel, Dmax, gamma, qmax, rat, Num,FALSE)
names(result)
result$mus
result$qqs
l <- length(gamma)
for(i in 1:l)print(result$MetaModel[[i]]$mu)
for(i in 1:l)print(result$MetaModel[[i]]$gamma)
for(i in 1:l)print(result$MetaModel[[i]]$'Meta-Model'$Nsupp)
```

A.9 SI_emp function

SI_emp *Function to calculate the empirical Sobol indices.*

Description

For each RKHS meta-model, this function calculates the empirical Sobol indices for all groups that are active in its support.

Usage

```
SI_emp(res,ErrPred)
```

Arguments

res List, includes a sequence of estimated meta-models, the solutions of the RKHS ridge group sparse or RKHS group lasso problems. It should have the same format as the output of one of the functions: [pen_MetMod](#), [RKHSMetMod](#) or [RKHSMetMod_qmax](#).

ErrPred Matrix or NULL. If matrix, each element of the matrix is the obtained prediction error associated with one RKHS meta-model in "res". It should have the same format as the output of the function [PredErr](#). Set as "NULL" by default.

Details

Details.

Value

If input `ErrPred` \neq "NULL", Vector of the empirical Sobol indices for the meta-model with the minimum Prediction error is returned. If `ErrPred` = "NULL", a list of the vectors is returned. Each vector is the obtained Sobol indices associated with one meta-model in "res".

Note:

Note.

Author(s)

Halaleh Kamari

References

References.

See Also

[PredErr](#), [pen_MetMod](#), [RKHSMetMod](#), [RKHSMetMod_qmax](#)

Examples

```
d <- 3
n <- 50;nT <- 50
library(lhs)
X <- maximinLHS(n, d);XT <- maximinLHS(nT, d)
c <- c(0.2,0.6,0.8)
F <- 1;for (a in 1:d) F <- F*(abs(4*X[,a]-2)+c[a])/(1+c[a])
FT <- 1;for (a in 1:d) FT <- FT*(abs(4*XT[,a]-2)+c[a])/(1+c[a])
sigma <- 0.2
epsilon <- rnorm(n,0,1);Y <- F + sigma*epsilon
epsilonT <- rnorm(nT,0,1);YT <- FT + sigma*epsilonT
Dmax <- 3
kernel <- "matern"
```

```
frc <- c(10)
gamma=c(.5,.01,.001)
res <- RKHSMetMod(Y,X,kernel,Dmax,gamma,frc,FALSE)
mu <- vector()
l <- length(gamma)
for(i in 1:length(frc))mu[i]=res[[i-1]*l+1]$mu
error <- PredErr(X,XT, YT,mu,gamma, res, kernel,Dmax)
SI.minErr <- SI_emp(res, error)
SI <- SI_emp(res, NULL)
```


Bibliography

Adamczak, R.

2005. Logarithmic sobolev inequalities and concentration of measure for convex functions and polynomial chaoses. *Bulletin of the Polish Academy of Sciences Mathematics*, 53. (Cited on pages 21, 51, 64, 82 and 83.)

Aronszajn, N.

1950. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3):337–404. (Cited on pages 8, 38 and 66.)

Bartlett, P. L., O. Bousquet, and S. Mendelson

2005. Local rademacher complexities. *Ann. Statist.*, 33(4):1497–1537. (Cited on pages 20, 50, 64 and 102.)

Bates, D. and D. Eddelbuettel

2013. Fast and elegant numerical linear algebra using the RcppEigen package. *Journal of Statistical Software*, 52(5):1–24. (Cited on page 130.)

Bednorz, W.

2014. Some remarks on the sudakov minoration. *ArXiv e-prints*. (Cited on page 79.)

Benoumechiara, N. and K. Elie-Dit-Cosaque

2019. Shapley effects for sensitivity analysis with dependent inputs: bootstrap and kriging-based algorithms. *ESAIM: Proceedings and Surveys*, 65:266–293. (Cited on pages 29 and 59.)

Berlinet, A. and C. Thomas-Agnan

2003. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Springer US. (Cited on pages 8, 38 and 66.)

Blatman, G. and B. Sudret

2011. Adaptive sparse polynomial chaos expansion based on least angle regression. *Journal of computational Physics*, 230:2345–2367. (Cited on pages 7, 10, 36, 39, 116 and 122.)

Bobkov, S. and M. Ledoux

1997. Poincaré’s inequalities and talagrand’s concentration phenomenon for the exponential distribution. *Probability Theory and Related Fields*, 107(3):383–400. (Cited on page 81.)

Borell, C.

1974. Convex measures on locally convex spaces. *Ark. Mat.*, 12(1-2):239–252. (Cited on pages 30, 59 and 79.)

- Borgonovo, E. and E. Plischke
2016. Sensitivity analysis: A review of recent advances. *European Journal of Operational Research*, 248(3):869 – 887. (Cited on pages 3 and 33.)
- Boucheron, S., G. Lugosi, and P. Massart
2000. A sharp concentration inequality with applications. *Random Struct. Algorithms*, 16:277–292. (Cited on page 104.)
- Boyd, S., N. Parikh, E. Chu, B. Peleato, and J. Eckstein
2011. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.*, 3(1):1–122. (Cited on pages 14, 26, 43, 56 and 127.)
- Breiman, L.
2001. Random forests. *Machine Learning*, 45(1):5–32. (Cited on pages 6 and 36.)
- Broto, B., F. Bachoc, M. Depecker, and J.-M. Martinez
2019. Sensitivity indices for independent groups of variables. *Mathematics and Computers in Simulation*, 163:19 – 31. (Cited on pages 29 and 59.)
- Bubeck, S.
2015. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3–4):231–357. (Cited on pages 14, 26, 43, 56 and 127.)
- Cacuci, D.
2003. *Sensitivity and Uncertainty Analysis, Theory*. New York: Chapman and Hall/CRC. (Cited on pages 3 and 33.)
- Cacuci, D. and Ionescu-Bujor, M. and I. Navon
2005. *Sensitivity and Uncertainty Analysis: Applications to Large-Scale Systems*, volume II. Boca Raton: CRC Press. (Cited on pages 3 and 33.)
- Chung, F. and L. Lu
2006. Concentration inequalities and martingale inequalities: a survey. *Internet Math.*, 3(1):79–127. (Cited on page 98.)
- de Rocquigny, E., N. Devictor, and S. Tarantola
2008. *Uncertainty in Industrial Practice: A Guide to Quantitative Uncertainty Management*. Wiley. (Cited on pages 3 and 33.)
- Dean, A. and S. Lewis
2006. *Screening: Methods for Experimentation in Industry, Drug Discovery, and Genetics*. Springer New York. (Cited on pages 3 and 33.)
- Durrande, N., D. Ginsbourger, and O. Roustant
2012. Additive covariance kernels for high-dimensional gaussian process modeling. *Annales de la faculté des sciences de Toulouse Mathématiques*, 21(3):481–499. (Cited on pages 7, 37 and 116.)

- Durrande, N., D. Ginsbourger, O. Roustant, and L. Carraro
2013. Anova kernels and rkhs of zero mean functions for model-based sensitivity analysis. *Journal of Multivariate Analysis*, 115:57 – 67. (Cited on pages 1, 7, 8, 9, 27, 29, 31, 37, 38, 57, 58, 62, 65, 114, 116, 118, 119, 121, 137 and 147.)
- Eddelbuettel, D.
2013. *Seamless R and C++ Integration with Rcpp*. Springer Publishing Company, Incorporated. (Cited on page 131.)
- Eddelbuettel, D. and R. Francois
2019. *RcppGSL: 'Rcpp' Integration for 'GNU GSL' Vectors and Matrices*. R package version 0.3.7. (Cited on page 130.)
- Faivre, R., B. Iooss, S. Mahévas, D. Makowski, and H. Monod
2013. *Analyse de sensibilité et exploration de modèles*, Collection Savoir-Faire. Editions Quae. (Cited on pages 3 and 33.)
- Fang, K.-T., R. Li, and A. Sudjianto
2005. *Design and Modeling for Computer Experiments (Computer Science & Data Analysis)*. Chapman & Hall/CRC. (Cited on pages 3 and 33.)
- Friedman, J. H.
1991. Multivariate adaptive regression splines. *Ann. Statist.*, 19(1):1–67. (Cited on pages 6, 12, 36 and 41.)
- Friedman, J. H.
2001. Greedy function approximation: A gradient boosting machine. *Ann. Statist.*, 29(5):1189–1232. (Cited on pages 6 and 36.)
- Fu, W. J.
1998. Penalized regressions: The bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416. (Cited on page 127.)
- Galassi, M. e. a.
2018. Gnu scientific library reference manual. (Cited on page 130.)
- Gentil, I., A. Guillin, and L. Miclo
2005. Modified logarithmic Sobolev inequalities and transportation inequalities. *Probability Theory and Related Fields*, 133 (3):409–436. (Cited on page 81.)
- Gentil, I., A. Guillin, and L. Miclo
2007. Modified logarithmic sobolev inequalities in null curvature. *Rev. Mat. Iberoamericana*, 23(1):235–258. (Cited on page 81.)
- Giraud, C.
2014. *Introduction to High-Dimensional Statistics*, Chapman & Hall/CRC Monographs on Statistics & Applied Probability. Taylor & Francis. (Cited on page 149.)

- Griewank, A. and A. Walther
2008. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, second edition. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics. (Cited on pages 3 and 33.)
- Gross, L.
1975. Logarithmic sobolev inequalities. *American Journal of Mathematics*, 97(4):1061–1083. (Cited on pages 21, 51, 64 and 81.)
- Guennebaud, G., B. Jacob, et al.
2010. Eigen v3. <http://eigen.tuxfamily.org>. (Cited on page 130.)
- Hastie, T., R. Tibshirani, and M. Wainwright
2015. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Chapman & Hall/CRC. (Cited on page 117.)
- Helton, J.
2008. *Uncertainty and Sensitivity Analysis for Models of Complex Systems.*, Pp. 207–228. Berlin, Heidelberg: Springer Berlin Heidelberg. (Cited on pages 3 and 33.)
- Hoeffding, W.
1948. A class of statistics with asymptotically normal distribution. *Ann. Math. Statist.*, 19(3):293–325. (Cited on page 62.)
- Homma, T. and A. Saltelli
1996. Importance measures in global sensitivity analysis of nonlinear models. *Reliability Engineering & System Safety*, 52(1):1 – 17. (Cited on pages 5, 35 and 115.)
- Huet, S. and M.-L. Taupin
2017. Metamodel construction for sensitivity analysis. *ESAIM: Procs*, 60:27–69. (Cited on pages 7, 12, 15, 18, 20, 27, 30, 37, 42, 45, 47, 49, 50, 57, 60, 62, 63, 64, 66, 69, 72, 77, 79, 114, 118, 147 and 153.)
- Iooss, B.
2011. Revue sur l’analyse de sensibilité globale de modèles numériques. *Journal de la Societe Française de Statistique*, 152(1):1–23. (Cited on pages 1 and 31.)
- Iooss, B. and P. Lemaître
2015. A review on global sensitivity analysis methods. In *Uncertainty management in Simulation-Optimization of Complex Systems: Algorithms and Applications*, C. Meloni and G. Dellino, eds. Springer. (Cited on pages 3 and 33.)
- Iooss, B. and C. Prieur
2019. Shapley effects for sensitivity analysis with correlated inputs: comparisons with Sobol’ indices, numerical estimation and applications. *International Journal for Uncertainty Quantification*, 9(5):493–514. (Cited on pages 29 and 59.)

- Kandasamy, K. and Y. Yu
2016. Additive approximations in high dimensional nonparametric regression via the salsa. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, Pp. 69–78. JMLR.org. (Cited on pages 12 and 41.)
- Kennedy, M. C. and A. O'Hagan
2001. Bayesian calibration of computer models. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 63(3):425–464. (Cited on pages 6 and 36.)
- Kimeldorf, G. S. and G. Wahba
1970. A correspondence between bayesian estimation on stochastic processes and smoothing by splines. *Ann. Math. Statist.*, 41(2):495–502. (Cited on pages 13, 43 and 120.)
- Kleijnen, J. P.
2009. Kriging metamodeling in simulation: A review. *European Journal of Operational Research*, 192(3):707 – 716. (Cited on pages 7, 37 and 116.)
- Kleijnen, J. P. C.
2007. *Design and Analysis of Simulation Experiments*, 1st edition. Springer Publishing Company, Incorporated. (Cited on pages 7, 37 and 116.)
- Koltchinskii, V. and M. Yuan
2010. Sparsity in multiple kernel learning. *Ann. Statist.*, 38(6):3660–3695. (Cited on pages 19, 20, 28, 29, 49, 50, 58, 63, 64 and 76.)
- Latała, R.
2014. Sudakov-type minoration for log-concave vectors. *Studia Mathematica*, 223(3):251–274. (Cited on pages 30, 59, 78 and 79.)
- Le Gratiet, L., C. Cannamela, and B. Iooss
2014. A bayesian approach for global sensitivity analysis of (multifidelity) computer codes. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):336–363. (Cited on pages 7, 37 and 116.)
- Le Gratiet, L., S. Marelli, and B. Sudret
2017. *Metamodel-Based Sensitivity Analysis: Polynomial Chaos Expansions and Gaussian Processes*, Pp. 1289–1325. Cham: Springer International Publishing. (Cited on pages 7, 37 and 116.)
- Ledoux, M.
1997. On talagrand's deviation inequalities for product measures. *ESAIM: Probability and Statistics*, 1:63–87. (Cited on page 81.)

Ledoux, M.

2001. *The Concentration of Measure Phenomenon*, Mathematical surveys and monographs. American Mathematical Society. (Cited on pages 20, 50 and 64.)

Ledoux, M. and M. Talagrand

1991. *Probability in Banach Spaces: isoperimetry and processes*. Berlin: Springer. (Cited on pages 81 and 102.)

Lin, Y. and H. H. Zhang

2006. Component selection and smoothing in multivariate nonparametric regression. *Ann. Statist.*, 34(5):2272–2297. (Cited on pages 12, 41 and 117.)

Liu, H., L. Wasserman, and J. D. Lafferty

2009. Nonparametric regression and classification with joint sparsity constraints. In *Advances in Neural Information Processing Systems 21*, D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, eds., Pp. 969–976. Curran Associates, Inc. (Cited on page 117.)

Marrel, A., B. Iooss, B. Laurent, and O. Roustant

2009. Calculations of sobol indices for the gaussian process metamodel. *Reliability Engineering & System Safety*, 94(3):742 – 751. (Cited on pages 7, 37 and 116.)

Massart, P.

2000. About the constants in talagrand’s concentration inequalities for empirical processes. *Ann. Probab.*, 28(2):863–884. (Cited on pages 20, 50 and 64.)

Massart, P. and J. Picard

2007. *Concentration Inequalities and Model Selection: Ecole d’Eté de Probabilités de Saint-Flour XXXIII - 2003*, Lecture Notes in Mathematics. Springer Berlin Heidelberg. (Cited on pages 74 and 104.)

Meier, L., S. van de Geer, and P. Bühlmann

2008. The group lasso for logistic regression. *Journal of the Royal Statistical Society. Series B*, 70(1):53–71. (Cited on page 127.)

Meier, L., S. van de Geer, and P. Bühlmann

2009. High-dimensional additive modeling. *Ann. Statist.*, 37(6B):3779–3821. (Cited on pages 11, 19, 41, 49, 63, 76 and 118.)

Mendelson, S.

2002. Geometric parameters of kernel machines. In *Computational learning theory (Sydney, 2002)*, volume 2375 of *Lecture Notes in Comput. Sci.*, Pp. 29–43. Springer, Berlin. (Cited on pages 16, 20, 46, 50, 64, 68 and 102.)

Milman, V. D. and G. Schechtman

1986. *Asymptotic Theory of Finite Dimensional Normed Spaces*. New York, NY, USA: Springer-Verlag New York, Inc. (Cited on page 84.)

- Oakley, J. E. and A. O'Hagan
2004. Probabilistic sensitivity analysis of complex models: a bayesian approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(3):751–769. (Cited on pages 6, 7, 36, 37 and 116.)
- Owen, A. B.
2014. Sobol' indices and shapley value. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):245–251. (Cited on pages 29 and 58.)
- Owen, A. B. and C. Prieur
2017. On shapley value for measuring importance of dependent inputs. *SIAM/ASA Journal on Uncertainty Quantification*, 5(1):986–1002. (Cited on pages 29 and 58.)
- Pisier, G.
1989. *The volume of convex bodies and Banach space geometry*, volume 94 of *Cambridge Tracts in Mathematics*. Cambridge University Press, Cambridge. (Cited on pages 20, 29, 50, 59, 64, 75, 77 and 78.)
- Raskutti, G., M. J. Wainwright, and B. Yu
2012. Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *J. Mach. Learn. Res.*, 13(1):389–427. (Cited on pages 11, 19, 20, 41, 49, 63, 64, 76, 77 and 117.)
- Raskutti, G., B. Yu, and M. J. Wainwright
2009. Lower bounds on minimax rates for nonparametric regression with additive sparsity and smoothness. In *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, eds., Pp. 1563–1570. Curran Associates, Inc. (Cited on page 117.)
- Ravikumar, P., J. Lafferty, H. Liu, and L. Wasserman
2009. Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030. (Cited on pages 11, 40 and 117.)
- Sacks, J., W. J. Welch, T. J. Mitchell, and H. P. Wynn
1989. Design and analysis of computer experiments. *Statist. Sci.*, 4(4):409–423. (Cited on pages 6 and 36.)
- Saitoh, S.
1988. *Theory of reproducing kernels and its applications*, Pitman research notes in mathematics series. Longman Scientific & Technical. (Cited on pages 8, 38 and 66.)
- Saltelli, A.
2002. Making best use of model evaluations to compute sensitivity indices. *Computer Physics Communications*, 145(2):280 – 297. (Cited on pages 5, 35 and 115.)

- Saltelli, A.
2008. *Global sensitivity analysis: the primer*. John Wiley. (Cited on pages 3 and 33.)
- Saltelli, A., K. Chan, and E. Scott
2009. *Sensitivity Analysis*. Wiley. (Cited on pages 1, 3, 31, 33 and 137.)
- Saltelli, A. and I. M. Sobol
1995. About the use of rank transformation in sensitivity analysis of model output. *Reliability Engineering & System Safety*, 50(3):225 – 239. (Cited on pages 3 and 33.)
- Saltelli, A. and S. Tarantola
2002. On the relative importance of input factors in mathematical models. *Journal of the American Statistical Association*, 97(459):702–709. (Cited on pages 4 and 33.)
- Schoutens, W.
2000. *Stochastic Processes and Orthogonal Polynomials*, Lecture Notes in Statistics. Springer New York. (Cited on pages 6, 36 and 115.)
- Shapley, L. S.
1953. *A value for n -person games*, P. 307–317. Princeton University Press, Princeton, NJ. (Cited on pages 29 and 58.)
- Shu, Y. and M. Strzelecki
2017. A characterization of a class of convex log-sobolev inequalities on the real line. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 54. (Cited on pages 21, 51, 64, 75, 82 and 83.)
- Sobol, I. M.
1993. Sensitivity estimates for nonlinear mathematical models. In *Sensitivity Estimates for Nonlinear Mathematical Models*. (Cited on pages 4, 5, 34, 35, 62, 114 and 115.)
- Sobol, I. M.
2001. Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Math. Comput. Simulation*, 55(1-3):271–280. The Second IMACS Seminar on Monte Carlo Methods (Varna, 1999). (Cited on pages 1 and 31.)
- Soize, C. and R. Ghanem
2004. Physical systems with random uncertainties: Chaos representations with arbitrary probability measure. *SIAM Journal on Scientific Computing*, 26(2):395–410. (Cited on pages 6, 36 and 115.)
- Song, E., B. L. Nelson, and J. Staum
2016. Shapley effects for global sensitivity analysis: Theory and computation.

- SIAM/ASA Journal on Uncertainty Quantification*, 4(1):1060–1083. (Cited on pages 29 and 58.)
- Storlie, C., H. Bondell, B. Reich, and H. Zhang
2011. Surface estimation, variable selection, and the nonparametric oracle property. *Statistica Sinica*, 21:679–705. (Cited on pages 6 and 36.)
- Storlie, C. B. and J. C. Helton
2008. Multiple predictor smoothing methods for sensitivity analysis: Description of techniques. *Reliability Engineering & System Safety*, 93(1):28 – 54. (Cited on pages 6 and 36.)
- Storlie, C. B., L. P. Swiler, J. C. Helton, and C. J. Sallaberry
2009. Implementation and evaluation of nonparametric regression procedures for sensitivity analysis of computationally demanding models. *Reliability Engineering & System Safety*, 94(11):1735 – 1763. (Cited on pages 6 and 36.)
- Sudret, B.
2008. Global sensitivity analysis using polynomial chaos expansions. *Reliability Engineering & System Safety*, 93(7):964 – 979. Bayesian Networks in Dependability. (Cited on pages 6, 36 and 115.)
- Talagrand, M.
1993. Regularity of infinitely divisible processes. *Ann. Probab.*, 21(1):362–432. (Cited on page 78.)
- Talagrand, M.
1994. The supremum of some canonical processes. *American Journal of Mathematics*, 116(2):283–325. (Cited on pages 20, 30, 50, 59, 64, 75, 79 and 80.)
- Touzani, S.
2011. *Response surface methods based on analysis of variance expansion for sensitivity analysis*. Theses, Université de Grenoble. (Cited on pages 6 and 36.)
- van de Geer, S., R. Gill, B. Ripley, S. Ross, B. Silverman, and M. Stein
2000. *Empirical Processes in M-Estimation*, Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press. (Cited on pages 20, 50, 64 and 76.)
- van der Vaart, A. W.
1998. *Asymptotic Statistics*, Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press. (Cited on pages 4, 34, 62 and 114.)
- Wahba, G.
1990. *Spline Models for Observational Data*. Society for Industrial and Applied Mathematics. (Cited on pages 12 and 41.)

Wahba, G., Y. Wang, C. gu, and B. Md

1995. Smoothing spline anova for exponential families, with application to the wisconsin epidemiological study of diabetic retinopathy. *Ann. Statist.*, 23. (Cited on pages 12 and 41.)

Welch, W. J., R. J. Buck, J. Sacks, H. P. Wynn, T. J. Mitchell, and M. D. Morris

1992. Screening, predicting, and computer experiments. *Technometrics*, 34(1):15–25. (Cited on pages 7, 37 and 116.)

Wiener, N.

1938. The homogeneous chaos. *American Journal of Mathematics*, 60(4):897–936. (Cited on pages 6, 36 and 115.)

Yang, Y. and H. Zou

2015. A fast unified algorithm for solving group-lasso penalize learning problems. *Statistics and Computing*, 25(6):1129–1141. (Cited on page 127.)

Yuan, M. and Y. Lin

2006. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67. (Cited on pages 23, 52, 121 and 127.)

Titre: Qualité prédictive des méta-modèles construits sur des espaces de Hilbert à noyau auto-reproduisant et analyse de sensibilité des modèles complexes

Mots clés: méta-modèle, des espaces de Hilbert à noyau auto-reproduisant, indice de Sobol, régression non-paramétrique, critère des moindres carrés pénalisé, majoration du risque

Résumé: Ce travail porte sur le problème de l'estimation d'un méta-modèle d'un modèle complexe, noté m . Le modèle m dépend de d variables d'entrée qui sont indépendantes et ont une loi connue. Le méta-modèle, noté f^* , approche la décomposition de Hoeffding de m et permet d'estimer ses indices de Sobol. Il appartient à un espace de Hilbert à noyau auto-reproduisant qui est construit comme une somme directe d'espaces de Hilbert (Durrande et al. (2013)). L'estimateur du f^* , noté \hat{f} , est calculé en minimisant un critère des moindres carrés pénalisé par la somme de la norme de Hilbert et de la norme empirique L^2 (Huet and Taupin (2017)). Ce travail se compose d'une partie théorique et d'une partie pratique. Dans la partie théorique, j'ai établi les majorations

du risque empirique L^2 et du risque quadratique de l'estimateur \hat{f} d'un modèle de régression où l'erreur est non-gaussienne et non-bornée. Dans la partie pratique, j'ai développé un package R appelé **RKHSMetaMod**, pour la mise en œuvre des méthodes d'estimation du méta-modèle f^* . Afin d'optimiser le temps de calcul et la mémoire de stockage, toutes les fonctions de ce package ont été écrites en utilisant les bibliothèques **GSL** et **Eigen** de C++ à l'exception d'une fonction qui est écrite en R. Elles sont ensuite interfacées avec l'environnement R afin de proposer un package facilement exploitable aux utilisateurs. La performance des fonctions du package en termes de qualité prédictive de l'estimateur et de l'estimation des indices de Sobol, est validée par une étude de simulation.

Title: Predictive quality of meta-models constructed on the reproducing kernel Hilbert spaces and sensitivity analysis of complex models

Keywords: meta-model, reproducing kernel Hilbert spaces, Sobol indices, non-parametric regression, penalized least-squares criterion, risk upper bound

Abstract: In this work, the problem of estimating a meta-model of a complex model, denoted m , is considered. The model m depends on d input variables that are independent and have a known law. The meta-model, denoted f^* , approximates the Hoeffding decomposition of m , and allows to estimate its Sobol indices. It belongs to a reproducing kernel Hilbert space (RKHS) which is constructed as a direct sum of Hilbert spaces (Durrande et al. (2013)). The estimator of f^* , denoted \hat{f} , is calculated by minimizing a least-squares criterion penalized by the sum of the Hilbert norm and the empirical L^2 -norm (Huet and Taupin (2017)). This work consists of a theoretical part and a practical part. In the theoretical part, I established upper bounds

of the empirical L^2 risk and the L^2 risk of the estimator \hat{f} in the regression framework with non-Gaussian and non-bounded error term. In the practical part, I developed an R package, called **RKHSMetaMod**, that implements the estimation methods of the meta-model f^* . In order to optimize the execution time and the storage memory, except for a function that is written in R, all of the functions of this package are written using C++ libraries **GSL** and **Eigen**. These functions are then interfaced with the R environment in order to propose an user friendly package. The performance of the package functions in terms of the predictive quality of the estimator and the estimation of the Sobol indices, is validated by a simulation study.

