

# ANALYSE DE SENSIBILITÉ EN DOMAINE CIRCULAIRE

Espéran Padonou<sup>1</sup> & Olivier Roustant<sup>1</sup>

<sup>1</sup> *Mines Saint-Etienne, UMR CNRS 6158, LIMOS, F-42023 Saint-Etienne, France.*  
*Email: padonou@emse.fr*

**Résumé** On considère le problème de reconstruction spatiale de données coûteuses sur le disque, motivé par des applications en environnement. Outre la spécificité du domaine d'étude, nous nous intéressons à l'interprétation et la visualisation qui sont des questions cruciales aux yeux des utilisateurs. En particulier, on souhaite décrire les variations horizontales, verticales, radiales et angulaires. Pour répondre à ces questions, nous utilisons la décomposition de Sobol-Hoeffding des processus gaussiens sur le disque avec une représentation en coordonnées cartésiennes ou polaires. Nous montrons comment l'utilisation des noyaux centrés permet de décomposer analytiquement la grandeur reconstruite en effets élémentaires dont l'influence est quantifiée à l'aide des indices de Sobol.

**Mots-clés.** Décomposition de Sobol-Hoeffding, processus gaussien, visualisation.

**Abstract.** Motivated by an application in environment, we consider Kriging models to recover data over the disk. In addition to the geometry of the domain, our focus is on interpretation and visualization, which represent two key points for users. In particular, the purpose is to describe horizontal, vertical, radial and angular variations. To address this problem, we use the Sobol-Hoeffding decomposition of Gaussian random field paths over the disk, represented either in Cartesian or polar coordinates. We show that the use of centred kernels allows to decompose analytically Kriging estimation into elementary terms whose importance are quantified with Sobol indices.

**Keywords.** Sobol-Hoeffding decomposition, Gaussian process, visualisation.

## 1 Motivation

Cette étude est motivée par une application en ingénierie de l'environnement, consistant à évaluer la quantité de méthane émise par une usine. Du fait de l'absence de capteurs adéquats pour mesurer la concentration de gaz, elle est simulée en fonction de la direction  $\theta$  (0 à 360 degrés) et de la vitesse  $\rho$  (0 à  $12m.s^{-1}$ ) du vent. Nous disposons de 242 simulations réalisées dans un scénario simplifié. Elles sont représentées sur le disque unité dont le bord correspond à la vitesse maximale (Figure 1). Fondé sur les modèles de la dynamique des fluides, le simulateur intègre le relief, les bâtiments, les cours d'eaux, les odeurs etc. En conditions réelles, le temps de calcul est donc important et le nombre de

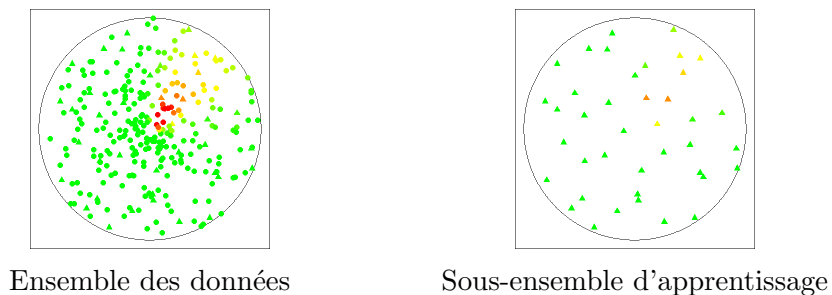


Figure 1: Données issues de 242 simulations de concentration de méthane

simulations limité. Pour reproduire cette difficulté, nos modèles seront estimés à partir d'un sous-ensemble de 30 points, les 212 autres servant d'ensemble test. La question est d'estimer la concentration de méthane aux points non explorés à partir du sous-ensemble d'apprentissage. Parmi les solutions envisagées, la régression par processus gaussiens, connue en géostatistique sous le nom krigeage, est retenue pour sa précision et la quantification d'incertitude spatiale qu'elle fournit (Rasmussen et Williams, 2006). Toutefois, le manque d'interprétabilité est une insuffisance que mentionnent souvent les utilisateurs.

Pour répondre à cette question, nous nous plaçons dans le cadre de la décomposition de Sobol-Hoeffding, un formalisme très utilisé en analyse de sensibilité des simulateurs coûteux. La décomposition de Sobol-Hoeffding donne les effets principaux et les termes d'interaction entre les entrées d'une fonction. Elle est fondée sur l'hypothèse d'entrées indépendantes, ce qui pose problème dans le cas du disque qui n'est pas un espace produit. Nous considérons alors deux possibilités pour réaliser l'analyse de sensibilité sur le disque:

- la première est de faire la décomposition sur le carré  $[-1, 1]^2$  en utilisant les coordonnées cartésiennes  $(x_1 = \rho \cos(\theta), x_2 = \rho \sin(\theta))$  et de la restreindre au disque;
- la seconde est d'utiliser l'espace produit des coordonnées polaires  $(x_1 = \rho, x_2 = \theta)$ . La difficulté est alors une modélisation adéquate de la variable directionnelle  $\theta$ .

L'apport original de cette étude est la formulation de la décomposition de Sobol-Hoeffding sur le disque, de façon analytique et non récursive. L'approche consiste à utiliser les processus gaussiens polaires introduits par Padonou et Roustant (2016), avec des noyaux de covariance d'intégrale nulle (Ginsbourger et al., 2015), correspondant à des trajectoires centrées. Cette propriété de centrage permet le calcul analytique des termes de la décomposition de Sobol-Hoeffding (Durrande 2013, Padonou 2016). Après avoir rappelé les conditions d'application de la décomposition de Sobol-Hoeffding, nous présenterons le modèle utilisé en expliquant en quoi il remplit ces conditions. La méthode sera ensuite illustrée à travers les données de la Figure 1, avec une visualisation des effets principaux (radial et angulaire) et d'interaction sur la concentration estimée.

## 2 Modèle

### 2.1 Décomposition de Sobol-Hoeffding

Soit  $D = \prod_{i=1}^d D_i$  avec  $D_i = [a_i, b_i]$  un hypercube de  $\mathbb{R}^d$ , et  $X_1, \dots, X_d$  des variables aléatoires indépendantes de  $D$ . Pour  $I \subseteq \{1 \dots d\}$  et  $\mathbf{x} = (x_1, \dots, x_d)$ , on désigne par  $\mathbf{x}_I$  le sous-vecteur de  $\mathbf{x}$  correspondant à  $I$ . Etant donné une fonction  $f \in L^2(D)$ ,  $f$  admet une décomposition unique dite décomposition de Sobol-Hoeffding (Sobol 1993):

$$f(\mathbf{x}) = f_0 + \sum_{i=1}^d f_i(x_i) + \sum_{i<j}^d f_{ij}(x_i, x_j) + \dots + f_{1\dots d}(x_1, \dots, x_d), \quad (1)$$

où  $f_0$  est une constante, et  $\forall I \subseteq \{1 \dots d\}$ ,  $f_I$  remplit la condition de centrage

$$\mathbb{E}(f_I(X_I)) = 0, \quad (2)$$

et la condition de non simplification

$$\mathbb{E}(f_{i_1 i_2}(X_{i_1}, X_{i_2}) | X_{i_1}) = \mathbb{E}(f_{i_1 i_2 i_3}(X_{i_1}, X_{i_2}, X_{i_3}) | X_{i_1}, X_{i_2}) = \dots = 0. \quad (3)$$

Les termes de (1) sont récursivement calculés via une procédure d'orthogonalisation:

$$\begin{aligned} f_0 &= \mathbb{E}(f(\mathbf{X})) \\ f_i(X_i) &= \mathbb{E}(f(\mathbf{X}) | X_i) - f_0, \quad i = 1 \dots d \\ f_{ij}(X_i, X_j) &= \mathbb{E}(f(\mathbf{X}) | X_i X_j) - f_i(X_i) - f_j(X_j) - f_0, \quad i, j = 1 \dots d, \quad \dots \end{aligned}$$

La condition de non simplification implique l'orthogonalité, d'où:

$$\text{var}(f(\mathbf{X})) = \sum_{i=1}^d \text{var}(f_i(X_i)) + \sum_{i<j}^d \text{var}(f_{ij}(X_i, X_j)) + \dots + \text{var}(f_{1\dots d}(X_1, \dots, X_d)).$$

La variance partielle  $V_I = \text{var}(f_I(X_I))$  quantifie la contribution de  $f_I$  dans la variance de  $f$ . En notant  $V = \text{var}(f(\mathbf{X}))$ , les ratios  $\frac{V_I}{V}$  sont appelés indices de Sobol.

### 2.2 Régression par processus gaussiens à trajectoires centrées

On suppose que la fonction  $f$  n'est évaluée qu'aux points d'apprentissage  $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$ , et on note  $\mathbf{y} = (Y_1, \dots, Y_n)$  le vecteur des réponses. Le krigeage est une méthode de lissage ou d'interpolation qui consiste à inférer les valeurs d'un processus gaussien conditionnellement aux observations. Nous considérons le modèle de krigeage suivant, inspiré de la décomposition de Sobol-Hoeffding des trajectoires des processus gaussiens (Ginsbourger et al., 2015), tronquée aux interactions d'ordre 2:

$$Z(\mathbf{x}) = \mu + Z_1(x_1) + \dots + Z_d(x_d) + \sum_{i<j} Z_{ij}(x_i, x_j) \quad (4)$$

où  $\mu$  est constant,  $Z_i(x_i) \sim GP(0, \sigma_i k_i)$  et  $Z_{ij}(x_i, x_j) \sim GP(0, \sigma_{i,j} k_i \otimes k_j)$  avec

1.  $Z_I$  et  $Z_J$  independants  $\forall I \neq J$ ,
2.  $\int_{D_i} k_i(u, v) d\nu_i(u) = 0 \quad \forall v \in D_i$

Défini par  $k(\mathbf{x}, \mathbf{x}') = \text{cov}(Z(\mathbf{x}), Z(\mathbf{x}'))$ , le noyau de  $Z$  joue un rôle clé. Dans (4), il est donné par  $k = \oplus_I k_I$ , mais on écrira  $k = \oplus k_I$  pour la simplicité des notations. Etant donné  $k$ , la prédiction en un nouveau point  $\mathbf{x}$  est alors donnée par l'espérance de  $Z(\mathbf{x})$  conditionnellement aux observations. Encore appelée moyenne de krigeage  $m(\mathbf{x})$ , cette valeur prédite est donnée par Rasmussen et Williams (2006):

$$m(\mathbf{x}) = \hat{\mu}(\mathbf{x}) + \mathbf{k}(\mathbf{x})^\top \mathbf{K}^{-1}(\mathbf{y} - \hat{\mu}(\mathbf{x})) \quad (5)$$

où  $\mathbf{K} = (k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))$  et  $\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}^{(i)}))$  sont respectivement la matrice de covariance aux points d'observation et le vecteur de covariance en  $\mathbf{x}$ ,  $i$  et  $j$  variant de 1 à  $n$ .

La condition 2 signifie que les noyaux  $k_i$  sont centrés. Lorsqu'elle est vérifiée, les trajectoires de  $Z$  sont aussi centrées et les termes de (4) s'interprètent de façon indépendante les uns des autres. Dans ce cas, la décomposition de Sobol-Hoeffding de la moyenne de krigeage est analytiquement donnée par Durrande (2013) et Padonou (2016):

$$m(\mathbf{x}) = m_0 + m_1(x_1) + \dots + m_d(x_d) + \sum_{i < j} m_{ij}(x_i, x_j) \quad (6)$$

où  $m_I(\mathbf{x}_I) = \alpha^\top k_I(\mathbf{x}_I)$ , avec  $\alpha = \mathbf{K}^{-1}(\mathbf{y} - \hat{\mu})$ . Les variances partielles sont données par:

$$\text{var}(m_I(\mathbf{x})) = \alpha^\top \mathbf{\Gamma}_I \alpha \quad (7)$$

avec  $\mathbf{\Gamma}_I = \odot_{i \in I} \mathbf{\Gamma}_i$ ,  $\odot$  désignant le produit terme à terme et  $\mathbf{\Gamma}_i = \int_{D_i} \mathbf{k}_i(x_i) \mathbf{k}_i(x_i)^\top d\nu_i(x_i)$ . De (6) et (7), il est possible de calculer sans récursivité les termes de la décomposition de Sobol-Hoeffding de la moyenne de krigeage et les indices de Sobol associés.

## 2.3 Exemples de noyaux centrés en dimension 1

Comparé aux modèles de krigeage usuels, la difficulté d'implémentation de (4) provient de la satisfaction de la condition de centrage des noyaux. Etant donné  $Z \sim GP(0, k)$  défini sur le segment  $[a, b]$  muni de la mesure  $\nu$ , une façon d'obtenir un noyau centré est d'utiliser le noyau du processus  $Z(x) - \int_a^b Z(x) d\nu(x)$  dont les trajectoires sont centrées (Ginsbourger et al., 2015). Son noyau  $k^*$  qui satisfait la condition de centrage est:

$$k^*(u, v) = k(u, v) - \int_{[a, b]} k(u, v) d\nu(u) - \int_{[a, b]} k(u, v) d\nu(v) + \iint_{[a, b]^2} k(u, v) d\nu(u) d\nu(v) \quad (8)$$

Dans le cas des variables directionnelles telles que l'angle  $\theta$  mentionné dans la section précédente (Figure 1), il est nécessaire d'utiliser un noyau  $2\pi$ -périodique. Soit  $T > 0$ ,

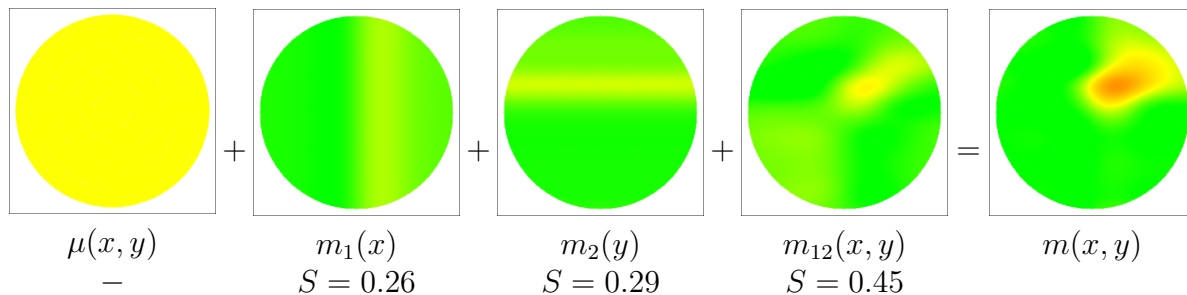
$\lambda$  la mesure uniforme sur  $[0, T]$  et  $\varphi : [0, \infty) \rightarrow \mathbb{R}$  une fonction  $T$ -périodique telle que  $k : (u, v) \mapsto \varphi(|u - v|)$  soit un noyau sur  $[0, T]^2$ . Alors, un noyau périodique centré sur  $[0, T]^2$  est donné par Padonou (2016):

$$k^*(u, v) = k(u, v) - s \quad (9)$$

où  $s = \int_0^T \varphi(t) d\lambda(t)$ . Des noyaux centrés sont ainsi obtenus par Padonou (2016) pour les variables directionnelles, à partir de la distance géodésique sur le cercle.

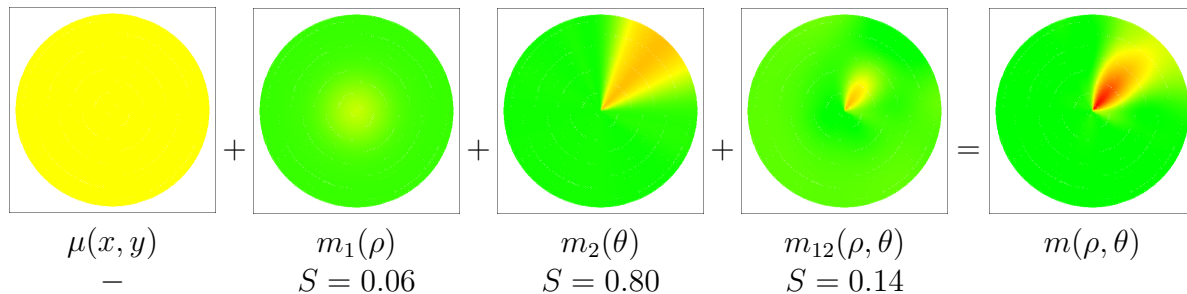
### 3 Application

On revient sur les données de la Figure 1. Les points sont représentés en coordonnées cartésiennes. A partir du sous-ensemble d'apprentissage, on estime par maximum de vraisemblance le modèle (4),  $k_1$  et  $k_2$  étant obtenus par centrage de la covariance  $C^2$  de Matérn (Equation (8) où  $\nu$  est la mesure uniforme sur  $[-1, 1]$ ). La décomposition qui en résulte (Equations (6) et (7)) est représentée graphiquement ci-dessous.



Décomposition de Sobol de la moyenne de krigeage (mesure uniforme sur  $[-1, 1]^2$ )

Cette décomposition montre qu'il est difficile de reconstruire la forme angulaire de la Figure 1 avec des processus horizontaux et verticaux. On considère ensuite les coordonnées polaires. On choisit la loi uniforme sur le disque, correspondant à une densité linéaire pour le rayon et la loi uniforme pour l'angle. Le noyau radial  $k_1$  est obtenu par centrage de la covariance  $C^2$  de Matérn. Le noyau angulaire  $k_2$  est donné par (9), à partir du noyau de Wendland (Padonou et Roustant, 2016). La décomposition qui en résulte est représentée ci-bas. Elle montre une meilleure reconstruction due à un effet angulaire prépondérant.



Décomposition de Sobol de la moyenne de krigeage (mesure uniforme sur le disque)

## 4 Discussion

Dans cette étude, nous avons utilisé la décomposition de Sobol-Hoeffding des processus gaussiens pour formuler un modèle de krigeage sur le disque. Fondée sur l'utilisation de noyaux centrés, l'approche permet d'interpréter la surface estimée en la décomposant en effets principaux et en termes d'interaction. L'application au cas des coordonnées polaires en montre l'efficacité. Cependant, la formulation proposée en coordonnées cartésiennes ne correspond pas à une décomposition sur le disque, mais sur le carré  $[-1, 1]^2$ , et pose la question suivante : comment réaliser une analyse de sensibilité d'une fonction  $f(x_1, x_2)$  définie sur  $[-1, 1]^2$  lorsque  $(x_1, x_2)$  suit la loi uniforme sur le disque ?

La difficulté provient de la forme du disque qui est incompatible avec l'hypothèse d'indépendance de la décomposition de Sobol-Hoeffding. Cette difficulté ne se résout pas simplement en considérant l'ANOVA généralisée au cas de variables dépendantes car la densité s'annule sur le complémentaire du disque sur le carré ("no hole condition", Owen et Prieur, 2016). Une solution est de quantifier l'importance des effets au moyen de la valeur de Shapley, au lieu des indices de Sobol, qui s'écrit explicitement en dimension 2 en fonction de la différence des indices de Sobol (Owen et Prieur, 2016). La visualisation graphique des effets principaux et des interactions reste possible, mais elle n'est plus associée à la décomposition de la variance.

## Bibliographie

- [1] N. Durrande, D. Ginsbourger, O. Roustant, L. Carraro, ANOVA kernels and RKHS of zero mean functions for model-based sensitivity analysis, *Journal of Multivariate Analysis*, 155, 57-67, 2013.
- [2] D. Ginsbourger, O. Roustant, D. Schuhmacher, N. Durrande, N. Lenz, On ANOVA decompositions of kernels and Gaussian random field paths in Monte Carlo and Quasi Monte Carlo Methods, *Proceedings in Mathematics & Statistics*, Springer, 163, 315-33, 2016.
- [3] E. Padonou, O. Roustant, Polar Gaussian processes and experimental designs in circular domains, *SIAM/ASA Journal on Uncertainty Quantification*, 4(1), 1014-1033, 2016.
- [4] E. Padonou, Statistical learning on circular domains for advanced process control in microelectronics, PhD thesis, Mines Saint-Etienne, 2016.
- [5] A. B. Owen, C. Prieur, On Shapley value for measuring importance of dependent inputs, preprint arXiv:1610.02080, 2016.
- [6] C.E. Rasmussen, C.K.I. Williams, Gaussian processes for machine learning, The MIT Press, 2006.
- [7] I. Sobol, Sensitivity estimates for non linear mathematical models, *Mathematical Modelling and Computational Experiments*, 1:407414, 1993.