



HAL
open science

Sensitivity analysis in general metric spaces

Fabrice Gamboa, Thierry Klein, Agnès Lagnoux, Leonardo Moreno

► **To cite this version:**

Fabrice Gamboa, Thierry Klein, Agnès Lagnoux, Leonardo Moreno. Sensitivity analysis in general metric spaces. 2020. hal-02044223v2

HAL Id: hal-02044223

<https://hal.science/hal-02044223v2>

Preprint submitted on 10 Feb 2020 (v2), last revised 19 Jan 2021 (v4)

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

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

Sensitivity analysis in general metric spaces

Fabrice Gamboa¹, Thierry Klein², Agnès Lagnoux³, and Leonardo Moreno⁴

¹Institut de Mathématiques de Toulouse; UMR5219. Université de Toulouse; CNRS. UT3, F-31062 Toulouse, France.

²Institut de Mathématiques de Toulouse; UMR5219. Université de Toulouse; ENAC - Ecole Nationale de l'Aviation Civile, Université de Toulouse, France

³Institut de Mathématiques de Toulouse; UMR5219. Université de Toulouse; CNRS. UT2J, F-31058 Toulouse, France.

⁴Departamento de Métodos Cuantitativos, FCEA, Universidad de la República, Uruguay

February 10, 2020

Abstract

In this paper, we introduce new indices adapted to outputs valued in general metric spaces. This new class of indices encompasses the classical ones; in particular, the so-called Sobol indices and the Cramér-von-Mises indices. Furthermore, we provide asymptotically Gaussian estimators of these indices based on U-statistics. Surprisingly, we prove the asymptotic normality straightforwardly. Finally, we illustrate this new procedure on a toy model and on two real-data examples.

Keywords: Sensitivity analysis, Cramér-von-Mises distance, Pick-Freeze method, U-statistics, general metric spaces.

1 Introduction

In the last decades, the use of computer code experiments to model physical phenomena has become a recurrent task for many applied researchers and engineers. In such simulations, it could be crucial to understand the global influence of one or several input variables on the output of the system. When considering these inputs as random elements, this problem is generally called (global) sensitivity analysis. We refer, for example to [5] or [19] for an overview on practical aspects of sensitivity analysis.

One of the most popular quantitative indicator to quantify the influence of some inputs is the so-called Sobol index. This index was first introduced in [20] and is well tailored, when the output space is \mathbb{R} . It compares thanks to the so-called Hoeffding decomposition (see [11]) the conditional variance of the output (knowing some of the input variables) with the global variance of the output. An efficient estimation of the Sobol indices can be performed through the Pick-Freeze method. For a description of this method and its theoretical study (consistency, central limit theorem, asymptotic efficiency, concentration inequalities and Berry-Esseen bounds), we refer to [12, 9] and references therein.

The case of vectorial outputs was first studied in [13] and tackled using principal component analysis of the output. In [8], the authors recover the indices proposed in [13] and showed that in some sense they are the only reasonable generalization of the classical Sobol indices in dimension greater than 2. Moreover, they provide the theoretical study of the Pick-Freeze estimators and extend their definitions to the case of outputs valued in a separable Hilbert space.

Since Sobol indices are based on the variance through the Hoeffding decomposition, they only quantify the input influence on the mean value of the computer code. Many authors proposed another way to compare the conditional distribution of the output knowing some of the inputs to the distribution of the output. In [16, 15], the authors considered higher moments to define new indices, whereas in [1, 2, 4], the use of divergences or distances between measures allows to define new indices. In [6], the authors used contrast functions to build goal-oriented indices. Although these works defined nice theoretical indices, the existence of an efficient statistical estimation procedure is still in most cases an open question. The case of vectorial-valued computer codes is considered in [10] where a sensitivity index based on the whole distribution of the output thanks to the Cramér-von-Mises distance is defined. The authors showed that the Pick-Freeze estimation procedure can be used providing an asymptotically Gaussian estimator of the index. This scheme requires $3N$ evaluations of the output code and leads to a convergence rate of order \sqrt{N} . This approach has been generalized in [7], where the authors considered computer codes valued in a compact Riemannian manifold. Once again, they used the Pick-Freeze scheme to provide a consistent estimator of their index requiring $4N$ evaluations of the output. Unfortunately, no central limit theorem was proved.

In this work, we build general indices for a code valued in a metric space and we provide asymptotically Gaussian estimator based on U-statistics requiring only $2N$ evaluations of the output code while keeping a convergence rate of \sqrt{N} . In addition, we explain that all the indices studied in [12, 9, 8, 10, 7] can be seen as particular cases of our framework. Hence, we improve the estimation scheme of [10] and [7] by reducing to $2N$ the number of evaluations of the code. Last but not least, thanks to the results of Hoeffding [11] on U-statistics, the asymptotic normality is proved straightforwardly.

The paper is organized as follows. Section 2 is dedicated to the definition of the new indices and the presentation of their estimation via U-statistics. In Section 3, we recover the classical indices used in sensitivity analysis. Furthermore, we extend the work of [7] and establish the central limit theorem that was not yet proved. We illustrate the procedure in Section 4 on a toy example and on two real-data models. The first application is about the Gaussian plume model and consists in quantifying the sensitivity of the contaminant concentration with respect to some input parameters. Second, an elliptical differential partial equation of type diffusive-advective transport is considered. In this setting, we proceed to the singular value decomposition of the solution and we perform a sensitivity analysis of the orthogonal matrix produced by the decomposition with respect to the equation parameters. Finally, some conclusions are given in Section 5.

2 General setting

We consider a regression function f (black-box code) defined on $E = E_1 \times E_2 \times \dots \times E_p$ and valued in a separable metric space (\mathcal{X}, d) . Here, $(E_1, \mathcal{A}_1), \dots, (E_p, \mathcal{A}_p)$ are measurable spaces. The output denoted by Z is given then by

$$Z = f(X_1, \dots, X_p), \quad (1)$$

where X_i is a random element of E_i and X_1, \dots, X_p are assumed to be mutually independent.

In [10], the authors studied for $\mathcal{X} = \mathbb{R}^k$ sensitivity indices of Z with respect to the inputs X_1, \dots, X_p based on its whole distribution (instead of considering only its second moment as done usually via the so-called Sobol indices). To do so, they introduced a family of test functions parametrized by a single index $t \in \mathbb{R}^k$ and defined by

$$T_t(Z) = \mathbb{1}_{\{Z \leq t\}},$$

where $\{Z \leq t\}$ means that $\{Z_1 \leq t_1, \dots, Z_k \leq t_k\}$.

Let \mathbf{u} be a subset of $I_p = \{1, \dots, p\}$ and let $\sim \mathbf{u}$ be its complementary in I_p ($\sim \mathbf{u} = I_p \setminus \{\mathbf{u}\}$). We define $X_{\mathbf{u}} = (X_i)_{i \in \mathbf{u}}$. Let also F be the distribution function of Z :

$$F(t) = \mathbb{P}(Z \leq t) = \mathbb{E}[\mathbb{1}_{\{Z \leq t\}}], \text{ for } t = (t_1, \dots, t_k) \in \mathbb{R}^k$$

and $F^{\mathbf{u}}$ be the conditional distribution function of Z conditionally on $X_{\mathbf{u}}$:

$$F^{\mathbf{u}}(t) = \mathbb{P}(Z \leq t | X_{\mathbf{u}}) = \mathbb{E}[\mathbb{1}_{\{Z \leq t\}} | X_{\mathbf{u}}], \text{ for } t = (t_1, \dots, t_k) \in \mathbb{R}^k.$$

Obviously, $\mathbb{E}[F^{\mathbf{u}}(t)] = F(t)$. Since for any fixed $t \in \mathbb{R}^k$, $T_t(Z)$ is a real-valued random variable, we can perform its Hoeffding decomposition:

$$\text{Var}(T_t(Z)) = F(t)(1 - F(t)) = \mathbb{E}[(F^{\mathbf{u}}(t) - F(t))^2] + \mathbb{E}[(F^{\sim \mathbf{u}}(t) - F(t))^2] + \text{Var}(R(t, \mathbf{u})) \quad (2)$$

where

$$R(t, \mathbf{u}) = T_t(Z) - \mathbb{E}[Y(t)T_t(Z) - (\mathbb{E}[T_t(Z)|X_{\mathbf{u}}] - \mathbb{E}[T_t(Z)]) - (\mathbb{E}[T_t(Z)|X_{\sim \mathbf{u}}] - \mathbb{E}[T_t(Z)])].$$

Then, the Cramér-von-Mises index is obtained by integrating in t with respect to the distribution of the output code Z :

$$S_{2,CVM}^{\mathbf{u}} = \frac{\int_{\mathbb{R}^k} \mathbb{E}[(F(t) - F^{\mathbf{u}}(t))^2] dF(t)}{\int_{\mathbb{R}^k} F(t)(1 - F(t)) dF(t)}. \quad (3)$$

In this example, the collection of the expectations $\mathbb{E}[T_t(Z)] = \mathbb{E}[\mathbb{1}_{\{Z \leq t\}}]$ ($t \in \mathbb{R}^k$) is parametrized by a single vectorial parameter t . Since its knowledge characterizes the distribution of Z , the previous indices depend as expected on the whole distribution of the output computer code. Using the Pick-Freeze methodology, the authors of [10] proposed an estimator which requires $3N$ evaluations of the output code leading to a convergence rate of \sqrt{N} .

This approach has been generalized in [7] to compact Riemannian manifolds replacing the indicator function of half-spaces $\mathbb{1}_{\{Z \leq t\}}$ parametrized by t by the indicator function of balls $\mathbb{1}_{\{Z \in B(a,b)\}}$ indexed by two parameters a and b . In their work, $B(a,b)$ stands for the ball of diameter \overline{ab} . In this last paper, a consistent estimation scheme based on $4N$ evaluations of the function is proposed. Nevertheless, the convergence rate of the estimator is not studied.

Now we aim at generalizing this methodology to any separable metric spaces and to any classes of test functions parametrized by a fixed number of elements of the metric space.

2.1 A new index

Generalizing the previous approach, we consider a family of test functions parametrized by m elements of \mathcal{X} with $m \in \mathbb{N}^*$. For any $a = (a_i)_{i=1, \dots, m} \in \mathcal{X}^m$, we consider the test functions

$$\begin{aligned} \mathcal{X}^m \times \mathcal{X} &\rightarrow \mathbb{R} \\ (a, x) &\mapsto T_a(x) \end{aligned}$$

We assume that $T_a(\cdot) \in L^2(\mathbb{P}^{\otimes m} \otimes \mathbb{P})$ where \mathbb{P} denotes the distribution of Z . Performing the Hoeffding decomposition on each test function $T_a(\cdot)$ and then integrating with respect to a using $\mathbb{P}^{\otimes m}$ leads to the definition of our new index.

Definition 2.1. The *general metric space sensitivity index* with respect to \mathbf{u} is defined by

$$S_{2,GMS}^{\mathbf{u}} := \frac{\int_{\mathcal{X}^m} \mathbb{E}_{X_{\mathbf{u}}} \left[(\mathbb{E}_Z[T_a(Z)] - \mathbb{E}_Z[T_a(Z)|X_{\mathbf{u}}])^2 \right] d\mathbb{P}^{\otimes m}(a)}{\int_{\mathcal{X}^m} \text{Var}(T_a(Z)) d\mathbb{P}^{\otimes m}(a)}, \quad (4)$$

where \mathbb{E}_U stands for the expectation with respect to the random variable U .

Proposition 2.2. *By construction, the new index lies in $[0, 1]$ and shares the same properties as the Sobol one:*

1. the different contributions sum to 1;
2. they are invariant by translation, by any isometry and by any non-degenerated scaling of the components of Z .

Particular examples

1. For $\mathcal{X} = \mathbb{R}$, $m = 0$ and T_a given by $T_a(x) = x$, one recovers the classical Sobol indices (see [21, 20]). For $\mathcal{X} = \mathbb{R}^k$ and $m = 0$, one can recover the index defined for vectorial outputs in [8, 13] by extending (4).
2. For $\mathcal{X} = \mathbb{R}^k$, $m = 1$ and T_a given by $T_a(x) = \mathbb{1}_{\{x \leq a\}}$, one recovers the index based on the Cramér-von-Mises distance defined in [10] and recalled in (3).
3. Consider that $\mathcal{X} = \mathcal{M}$ is a manifold, $m = 2$ and T_a is given by $T_a(x) = \mathbb{1}_{\{x \in B(a_1, a_2)\}}$, where $B(a_1, a_2)$ will stand for the ball of diameter $\overline{a_1 a_2}$. Here, one recovers the index defined in [7].

2.2 Estimation procedure via U-statistics

Following the so-called Pick-Freeze scheme, let $X^{\mathbf{u}}$ be the random vector such that $X_{\mathbf{u}}^{\mathbf{u}} = X_{\mathbf{u}}$ and $X_i^{\mathbf{u}} = X'_i$ if $i \notin \mathbf{u}$ where X'_i is an independent copy of X_i . Then, setting

$$Z^{\mathbf{u}} := f(X^{\mathbf{u}}), \quad (5)$$

a classical computation leads to the following relationship (see, e.g., [12]):

$$\text{Var}(\mathbb{E}[T_a(Z)|X_{\mathbf{u}}]) = \text{Cov}(T_a(Z), T_a(Z^{\mathbf{u}})).$$

Let us define $\mathbf{Z} = (Z, Z^{\mathbf{u}})^{\top}$ and consider $(m+2)$ i.i.d. copies of \mathbf{Z} denoted by $(\mathbf{Z}_i, i = 1, \dots, m+2)$. In the sequel, $\mathbb{P}_{\mathbf{Z}}^{\mathbf{u}}$ stands for the law of $\mathbf{Z} = (Z, Z^{\mathbf{u}})^{\top}$. Then the integrand in the numerator of (4) rewrites as

$$\begin{aligned} \mathbb{E} \left[(\mathbb{E}[T_a(Z)] - \mathbb{E}[T_a(Z)|X_{\mathbf{u}}])^2 \right] &= \mathbb{E}_{Z_1, \dots, Z_m} [\text{Var}(\mathbb{E}[T_a(Z_{m+1})|X_{\mathbf{u}}])] \\ &= \mathbb{E}_{Z_1, \dots, Z_m} [\text{Cov}_{\mathbf{Z}_{m+1}}(T_{Z_1, \dots, Z_m}(Z_{m+1}), T_{Z_1, \dots, Z_m}(Z_{m+1}^{\mathbf{u}}))]. \end{aligned}$$

Here the notation \mathbb{E}_Z (resp. Cov_Z) stands for the expectation (resp. the covariance) with respect to the law of the random variable Z .

Now, for any $1 \leq i \leq m+2$, we let $\mathbf{z}_i = (z_i, z_i^{\mathbf{u}})$ and we define

$$\begin{aligned}\Phi_1(\mathbf{z}_1, \dots, \mathbf{z}_{m+1}) &:= T_{z_1, \dots, z_m}(z_{m+1})T_{z_1, \dots, z_m}(z_{m+1}^{\mathbf{u}}) \\ \Phi_2(\mathbf{z}_1, \dots, \mathbf{z}_{m+2}) &:= T_{z_1, \dots, z_m}(z_{m+1})T_{z_1, \dots, z_m}(z_{m+2}^{\mathbf{u}}) \\ \Phi_3(\mathbf{z}_1, \dots, \mathbf{z}_{m+1}) &:= T_{z_1, \dots, z_m}(z_{m+1})^2 \\ \Phi_4(\mathbf{z}_1, \dots, \mathbf{z}_{m+2}) &:= T_{z_1, \dots, z_m}(z_{m+1})T_{z_1, \dots, z_m}(z_{m+2}).\end{aligned}$$

We further set

$$m(1) = m(3) = m+1 \quad \text{and} \quad m(2) = m(4) = m+2 \quad (6)$$

and we define, for $j = 1, \dots, 4$,

$$I(\Phi_j) := \int_{\mathcal{X}^{m(j)}} \Phi_j(\mathbf{z}_1, \dots, \mathbf{z}_{m(j)}) d\mathbb{P}_2^{u, \otimes m(j)}(\mathbf{z}_1, \dots, \mathbf{z}_{m(j)}). \quad (7)$$

Finally, we introduce the application Ψ from \mathbb{R}^4 to \mathbb{R} defined by

$$\begin{aligned}\Psi : \quad \mathbb{R}^4 &\rightarrow \mathbb{R} \\ (x, y, z, t) &\mapsto \frac{x-y}{z-t}.\end{aligned} \quad (8)$$

Then, $S_{2,GMS}^{\mathbf{u}}$ can be rewritten as

$$S_{2,GMS}^{\mathbf{u}} = \Psi(I(\Phi_1), I(\Phi_2), I(\Phi_3), I(\Phi_4)). \quad (9)$$

The previous expression of $S_{2,GMS}^{\mathbf{u}}$ will allow to perform easily its estimation. Following Hoeffding [11], we replace the functions Φ_1, Φ_2, Φ_3 and Φ_4 by their symmetrized version $\Phi_1^s, \Phi_2^s, \Phi_3^s$ and Φ_4^s :

$$\Phi_j^s(\mathbf{z}_1, \dots, \mathbf{z}_{m(j)}) = \frac{1}{(m(j))!} \sum_{\tau \in \mathcal{S}_{m(j)}} \Phi_j(\mathbf{z}_{\tau(1)}, \dots, \mathbf{z}_{\tau(m(j))})$$

for $j = 1, \dots, 4$ where \mathcal{S}_k is the symmetric group of order k (that is the set of all permutations on I_k). For $j = 1, \dots, 4$, the integrals $I(\Phi_j^s)$ are naturally estimated by U-statistics of order $m(j)$. More precisely, we consider an i.i.d. sample $(\mathbf{Z}_1, \dots, \mathbf{Z}_N)$ ($N \geq 1$) with distribution $\mathbb{P}_2^{\mathbf{u}}$ and, for $j = 1, \dots, 4$, we define the U-statistics

$$U_{j,N} := \binom{N}{m(j)}^{-1} \sum_{1 \leq i_1 < \dots < i_{m(j)} \leq N} \Phi_j^s(\mathbf{Z}_{i_1}, \dots, \mathbf{Z}_{i_{m(j)}}). \quad (10)$$

Theorem 7.1 in [11] ensures that $U_{j,N}$ converges in probability to $I(\Phi_j)$ for any $j = 1, \dots, 4$. Moreover, one may also prove that the convergence holds almost surely proceeding as in the proof of Lemma 6.1 in [10]. Then we estimate $S_{2,GMS}^{\mathbf{u}}$ by

$$\widehat{S}_{2,GMS}^{\mathbf{u}} := \frac{U_{1,N} - U_{2,N}}{U_{3,N} - U_{4,N}} = \Psi(U_{1,N}, U_{2,N}, U_{3,N}, U_{4,N}). \quad (11)$$

Our main result follows.

Theorem 2.3. *If for $j = 1, \dots, 4$, $\mathbb{E} \left[\Phi_j^s(\mathbf{Z}_1, \dots, \mathbf{Z}_{m(j)})^2 \right] < \infty$ then*

$$\sqrt{N} \left(\widehat{S}_{2,GMS}^{\mathbf{u}} - S_{2,GMS}^{\mathbf{u}} \right) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}_1(0, \sigma^2) \quad (12)$$

where the asymptotic variance σ^2 is given by (13) in the proof below.

Proof of Theorem 2.3. The first step of the proof is to apply Theorem 7.1 of [11] to the random vector $(U_{1,N}, U_{2,N}, U_{3,N}, U_{4,N})^\top$. By Theorem 7.1 and Equations (6.1)-(6.3) in [11], it follows that

$$\sqrt{N} \left(\begin{pmatrix} U_{1,N} \\ U_{2,N} \\ U_{3,N} \\ U_{4,N} \end{pmatrix} - \begin{pmatrix} I(\Phi_1^s) \\ I(\Phi_2^s) \\ I(\Phi_3^s) \\ I(\Phi_4^s) \end{pmatrix} \right) \xrightarrow[N \rightarrow \infty]{\mathcal{L}} \mathcal{N}_4(0, \Gamma)$$

where Γ is the square matrix of size 4 given by

$$\Gamma(i, j) := m(i)m(j)\text{Cov}(\mathbb{E}[\Phi_i^s(\mathbf{Z}_1, \dots, \mathbf{Z}_{m(i)})|\mathbf{Z}_1], \mathbb{E}[\Phi_j^s(\mathbf{Z}_1, \dots, \mathbf{Z}_{m(j)})|\mathbf{Z}_1]).$$

Now, it remains to apply the so-called Delta method (see [25]) with the function Ψ defined by (8). Thus, one gets the asymptotic behavior in Theorem 2.3 where σ^2 is given by

$$\sigma^2 := g^\top \Gamma g \tag{13}$$

with $g = \nabla \Psi(I(\Phi_1^s), I(\Phi_2^s), I(\Phi_3^s), I(\Phi_4^s))$ and $\nabla \Psi = (z-t)^{-2} (z-t, -z+t, -x+y, x-y)^\top$. \square

Notice that we consider $(m+2)$ copies of \mathbf{Z} in the definition of $S_{2,GMS}^u$ (see (9)). Nevertheless, the estimation procedure only requires a N sample of \mathbf{Z} (see (11)) that means only $2N$ evaluations of the black-box code which constitutes an appealing advantage of the method presented in this paper. Moreover, the required number of calls to the black-box code is independent of the size m of the class of tests functions unlike in [10] or in [7] where $(m+2) \times N$ calls of the computer code were necessary. In addition, the proof of the asymptotic normality in Theorem 2.3 is elementary and does not rely anymore on the use of the sophisticated fonctionnal Delta method as in [10].

2.3 Comments

Considering an output code f , one may consider different choices of the family $(T_a)_{a \in \mathcal{X}^m}$ of functions indexed by $a \in \mathcal{X}^m$ leading to very different indices. The choice of the family must be induced by the aim of the practitioner. To quantify the output sensitivity around the mean, one should consider the classical Sobol indices based on the variance and corresponding to the first example of Section 2.1. Otherwise, interested in the sensitivity of the whole distribution, one should prefer a family of functions that characterizes the distribution. For instance, in the previous second example of Section 2.1, the functions T_a are the indicator functions of half-lines and yield the Cramér-von-Mises indices.

Moreover, since in the estimation procedure the number of output calls is independent of the choice of the family $(T_a)_{a \in \mathcal{X}^m}$, one can consider and estimate simultaneously several indices with no-extra cost. In fact, the only computational challenge relies in our capability to evaluate the functions Φ on the sample.

3 Applications in classical frameworks and beyond

3.1 Particular cases

Sobol indices For $\mathcal{X} = \mathbb{R}$, $m = 0$ and the test functions T_a given by $T_a(x) = x$, our method provides a new estimator based on U-statistics for the classical Sobol index. In that case, the

estimator is given by (11) and, for $j = 1, \dots, 4$, the $U_{j,N}$'s are given by

$$\begin{aligned} U_{1,N} &= \frac{1}{N} \sum_{i=1}^N Z_i Z_i^{\mathbf{u}} \\ U_{2,N} &= \frac{1}{N(N-1)} \left(\sum_{i=1}^N Z_i \sum_{i=1}^N Z_i^{\mathbf{u}} - \sum_{i=1}^N Z_i Z_i^{\mathbf{u}} \right) =: \frac{1}{N(N-1)} (\tilde{U}_{2,N} - \tilde{V}_{2,N}) \\ U_{3,N} &= \frac{1}{N} \sum_{i=1}^N Z_i^2 \\ U_{4,N} &= \frac{1}{N(N-1)} \left(\left(\sum_{i=1}^N Z_i \right)^2 - \sum_{i=1}^N Z_i^2 \right) =: \frac{1}{N(N-1)} (\tilde{U}_{4,N} - \tilde{V}_{4,N}) \end{aligned}$$

leading to

$$\hat{S}_{2,GMS}^{\mathbf{u}} = \frac{U_{1,N} - U_{2,N}}{U_{3,N} - U_{4,N}} = \Psi(U_{1,N}, U_{2,N}, U_{3,N}, U_{4,N})$$

while in [9], the classical Pick-Freeze estimator $\hat{S}^{\mathbf{u}}$ of $S_{2,GMS}^{\mathbf{u}}$ is given by

$$\hat{S}^{\mathbf{u}} = \frac{U_{1,N} - (1 - 1/N^2)\tilde{U}_{2,N}}{U_{3,N} - (1 - 1/N^2)\tilde{U}_{4,N}} = \Psi(U_{1,N}, (1 - 1/N^2)\tilde{U}_{2,N}, U_{3,N}, (1 - 1/N^2)\tilde{U}_{4,N}) \quad (14)$$

and takes into account the diagonal terms. Both procedures require $2N$ evaluations of the black-box code and have the same rate of convergence. The estimators are slightly different which induces different asymptotic variances. Finally, one may improve the procedures using the information of the whole sample leading to the analog version of the estimation $\hat{T}^{\mathbf{u}}$ given in [9, Eq.(6)]:

$$\hat{T}^{\mathbf{u}} = \frac{\frac{1}{N} \sum_{i=1}^N Y_i Y_i^{\mathbf{u}} - \left(\frac{1}{N} \sum_{j=1}^N \frac{Y_j + Y_j^{\mathbf{u}}}{2} \right)^2}{\frac{1}{N} \sum_{i=1}^N \frac{(Y_i)^2 + (Y_i^{\mathbf{u}})^2}{2} - \left(\frac{1}{N} \sum_{i=1}^N \frac{Y_i + Y_i^{\mathbf{u}}}{2} \right)^2}. \quad (15)$$

The sequence of estimators $\hat{T}^{\mathbf{u}}$ is asymptotically efficient in the Cramér-Rao sense (see [9, Proposition 2.5]). In this paper, we also could have constructed a new estimator $\hat{T}_{2,GMS}^{\mathbf{u}}$ analog version of $\hat{S}_{2,GMS}^{\mathbf{u}}$ taking into account the whole information contained in the sample. Anyway, based on the same initial design as $\hat{S}^{\mathbf{u}}$ and $\hat{T}^{\mathbf{u}}$, neither $\hat{S}_{2,GMS}^{\mathbf{u}}$ nor $\hat{T}_{2,GMS}^{\mathbf{u}}$ will be asymptotically efficient. Nevertheless, the estimation procedure proposed in this paper outperforms the procedure presented in [10, 7] as soon as $m \geq 1$.

Sobol indices for multivariate outputs For $\mathcal{X} = \mathbb{R}^k$ and $m = 0$, one may realize the same analogy between the estimation procedure proposed in this paper and that in [8].

Cramér-von-Mises indices For $\mathcal{X} = \mathbb{R}^k$, $m = 1$ and the test functions T_a given by $T_a(x) = \mathbb{1}_{\{x \leq a\}}$, we outperform the central limit theorem proved in [10]. Indeed, the estimator proposed in [10] requires $3N$ evaluations of the computer code versus only $2N$ in our new procedure. In addition, the proof therein is based on the powerful but complex functional Delta method while the proof of Theorem 2.3 is an elementary application of Theorem 7.1 in [11] combined with the classical Delta method.

3.2 Compact manifolds

A particular framework is the case when the output space is a compact Riemannian manifold \mathcal{M} . In [7], a similar index to $S_{2,GMS}^{\mathbf{u}}$ is studied in this special context, taking $T_a(x) = \mathbb{1}_{\{x \in B(a_1, a_2)\}}$ as test functions. The authors showed that, under some restrictions on the underlying probability measure and the Riemannian manifold, the family of balls $(B(a_1, a_2))_{(a_1, a_2) \in \mathcal{M}}$ is a determining class, that is, if two probability measures P_1 and P_2 on \mathcal{M} coincide on all the events of this family, then $P_1 = P_2$. By this property, they proved that if their index, denoted $B_2^{\mathbf{u}}$, vanishes then the distributions of $T_a(Z)$ and $T_a(Z)|X_{\mathbf{u}}$ coincide. Further, in the last paper, the performance of $B_2^{\mathbf{u}}$ in Riemannian manifolds immersed in \mathbb{R}^d with $d = 2, 3$ and on the cone of positive definite matrices (manifold) is analyzed. Last, an exponential inequality for the estimator $\hat{B}_2^{\mathbf{u}}$ of $B_2^{\mathbf{u}}$ is provided together with the almost sure convergence that is deduced from. Unfortunately, no central limit theorem is given.

As a particular case of $S_{2,GMS}^{\mathbf{u}}$, the asymptotic distribution of $\hat{B}_2^{\mathbf{u}}$ can be found from Theorem 2.3. Given x , since $(a_1, a_2) \mapsto T_{(a_1, a_2)}(x)$ is a symmetric function and $m = 2$, it is verified that,

$$\begin{aligned}\Phi_1(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) &= \mathbb{1}_{\{\mathbf{z}_3, \mathbf{z}_3^{\mathbf{u}} \in B(\mathbf{z}_1, \mathbf{z}_2)\}}, \\ \Phi_2(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3, \mathbf{z}_4) &= \mathbb{1}_{\{\mathbf{z}_3, \mathbf{z}_4^{\mathbf{u}} \in B(\mathbf{z}_1, \mathbf{z}_2)\}}, \\ \Phi_3(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) &= \mathbb{1}_{\{\mathbf{z}_3 \in B(\mathbf{z}_1, \mathbf{z}_2)\}}, \\ \Phi_4(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3, \mathbf{z}_4) &= \mathbb{1}_{\{\mathbf{z}_3, \mathbf{z}_4 \in B(\mathbf{z}_1, \mathbf{z}_2)\}}.\end{aligned}$$

In this setting, the limiting covariance matrix Γ is given by $\Gamma(i, j) = m(i)m(j)\text{Cov}(L_i, L_j)$, for $i, j = 1, \dots, 4$ where

$$\begin{aligned}L_1 &= \frac{1}{6} \sum_{\tau \in \mathcal{S}_3} \mathbb{P}(Z_{\tau_3}, \mathbf{z}_{\tau_3}^{\mathbf{u}} \in B(Z_{\tau_1}, Z_{\tau_2})|Z_1), \\ L_2 &= \frac{1}{24} \sum_{\tau \in \mathcal{S}_4} \mathbb{P}(Z_{\tau_3}, Z_{\tau_4}^{\mathbf{u}} \in B(Z_{\tau_1}, Z_{\tau_2})|Z_1), \\ L_3 &= \frac{1}{6} \sum_{\tau \in \mathcal{S}_3} \mathbb{P}(Z_{\tau_3} \in B(Z_{\tau_1}, Z_{\tau_2})|Z_1), \\ L_4 &= \frac{1}{24} \sum_{\tau \in \mathcal{S}_4} \mathbb{P}(Z_{\tau_3}, Z_{\tau_4} \in B(Z_{\tau_1}, Z_{\tau_2})|Z_1).\end{aligned}$$

It is also possible to take other test functions T_a (with $m = 2$). In this case, the index can be built in a more general metric spaces (\mathcal{M}, d) . For instance, $T_a(x) = \mathbb{1}_{\{x \in B(a_1, a_2) \cup B(a_2, a_1)\}}$ or $T_a(x) = \mathbb{1}_{\{x \in B(a_1, a_2) \cap B(a_2, a_1)\}}$.

4 Numerical applications

4.1 A non linear model

In this section, we illustrate and compare the different estimation procedures based on the Pick-Freeze scheme and the U-statistics for the classical the Sobol indices on the following toy model:

$$Z = \exp\{X_1 + 2X_2\}, \tag{16}$$

where X_1 and X_2 are independent standard Gaussian random variables. The distribution of Z is log-normal and we can derive both its density and its distribution functions:

$$f_Z(y) = \frac{1}{\sqrt{10\pi y}} e^{-(\ln y)^2/10} \mathbb{1}_{\mathbb{R}^+}(y) \quad \text{and} \quad F_Y(y) = \Phi\left(\frac{\ln y}{\sqrt{5}}\right).$$

Here, Φ stands for the distribution function of the standard Gaussian random variable. We have $p = 2$ and tedious exact computations lead to closed forms of the Sobol indices S^1 and S^2 :

$$S^1 = \frac{1 - e^{-1}}{e^4 - 1} \approx 0.0118 \quad \text{and} \quad S^2 = \frac{e^3 - e^{-3}}{e^4 - 1} \approx 0.3738.$$

Further, the Cramér-von-Mises indices $S_{2,CVM}^1$ and $S_{2,CVM}^2$ are also explicitly computable:

$$S_{2,CVM}^1 = \frac{6}{\pi} \arctan 2 - 2 \approx 0.1145 \quad \text{and} \quad S_{2,CVM}^2 = \frac{6}{\pi} \arctan \sqrt{19} - 2 \approx 0.5693.$$

The reader is referred to [10] for the details of these computations.

In Figure 1, we compare the estimations of the two first order Sobol indices and the estimations of the two first order Cramér-von-Mises indices obtained by both estimation procedures (U-statistics and Pick-Freeze). The total number of calls of the computer code range from $n = 100$ to 500000. To have a fair comparison, when estimating the Sobol indices with both methodologies (U-statistics and Pick-Freeze) and the Cramér-von-Mises indices using U-statistics, we have considered samples of size $N = n/(p + 1)$. In contrast, when estimating the Cramér-von-Mises indices using the Pick-Freeze scheme, we have considered samples of size $N = n/(p + 2)$. Then, each estimation requires a total number n of evaluations of the code. We observe that both methods converge and give precise results for large sample sizes. In addition, the estimation procedure with U-statistics outperforms the Pick-Freeze one as soon as $m \geq 1$ as expected. Such a better performance increases with the number m of parameters of the tests functions family.

4.2 The Gaussian plume model

In this section, the model under study concerns about a point source that emits contaminant into a uni-directional wind in an infinite domain. Such a model is also applied, for instance, to volcanic eruptions, pollen and insect dispersals, and is called the Gaussian plume model (GPM) (see, e.g., [3, 24]). The GPM assumes that atmospheric turbulence is stationary and homogeneous. Naturally, in Earth Sciences, it is crucial to analyze the sensitivity of the output of the GPM model regarding the input parameters (see [14, 17]).

The model parameters are represented in Figure 2. The contaminant is emitted at a constant rate Q and the wind direction is denoted by $\mathbf{u} = (u, 0, 0)$ (with $u \geq 0$) while the effective height is $H = h(1 + \delta)$ where h is the stack height and δh is the plume rise.

Then the contaminant concentration at location (x, y, z) is given by

$$C(x, y, z) = \frac{Q}{4\pi ur(x)} e^{-\frac{y^2}{4r(x)}} \left(e^{-\frac{(z-H)^2}{4r(x)}} + e^{-\frac{(z+H)^2}{4r(x)}} \right),$$

where r is a parametric function given by $r(x) = \frac{1}{u} \int_0^x K(v) dv$, the function K being the eddy diffusion. In this section, we investigate the particular two-dimensional case: the height is considered as zero (at ground level). In addition, we suppose that $r(x) = Kx/u$ where K is a constant. Hence, the contaminant concentration at location $(x, y, 0)$ rewrites as:

$$C(x, y, 0) = \frac{Q}{2\pi Kx} e^{-\frac{u(y^2+H^2)}{4Kx}}. \quad (17)$$

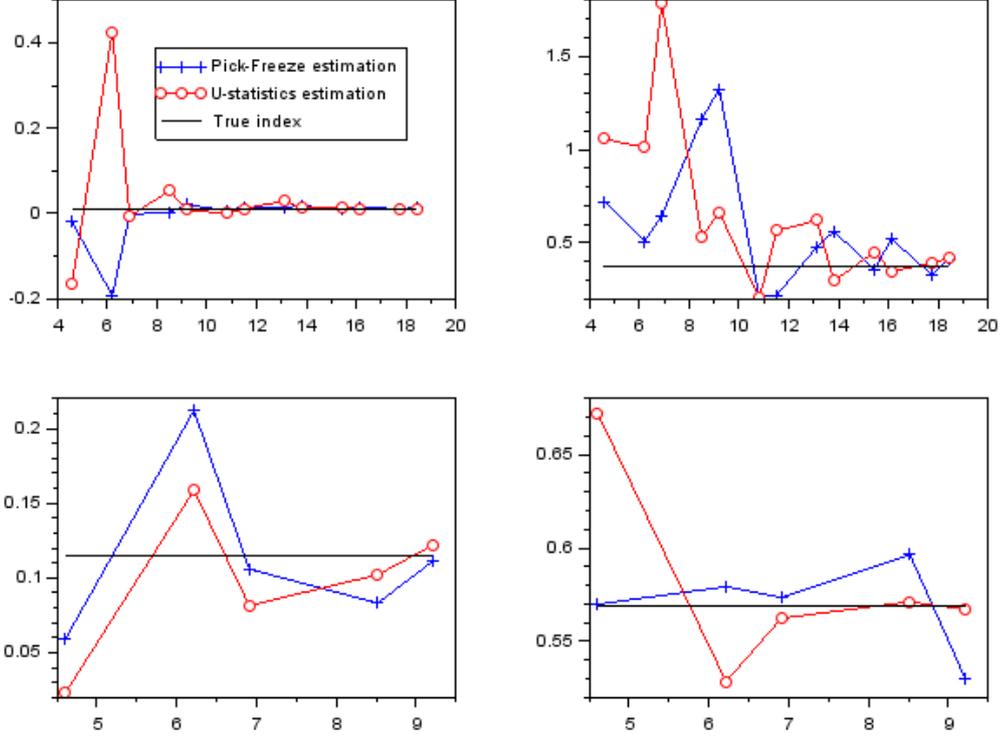


Figure 1: Non-linear model (16). Convergence of both methods when the total number of calls of the computer code increases. The two first order Sobol indices have been represented from left to right at the top row while the two first order Cramér-von-Mises indices have been represented from left to right at the bottom row. Several total number of calls of the computer code have been considered: $N = 100, 500, 1000, 5000, 10000, 50000, 100000,$ and 500000 . When estimating the Sobol indices with both methodologies (U-statistics and Pick-Freeze) and the Cramér-von-Mises indices with the U-statistics, we have considered samples of size $N = n/3$. In contrast, when estimating the Cramér-von-Mises indices using the Pick-Freeze scheme, we have considered samples of size $N = n/4$. The x -axis is in logarithmic scale.

Now we wish to perform a sensitivity analysis on the contaminant concentration with respect to the uncertain inputs Q , K , and u , while the altitude plume parameter H is fixed in advance. In this setting, the function f that defines the output of interest in (1) is then given by:

$$f: \quad \mathbb{R}^3 \quad \rightarrow \quad L^2(\mathbb{R}^2) \\ (Q, K, u) \quad \mapsto \quad f(Q, K, u) = (C(x, y, 0))_{(x, y) \in \mathbb{R}^2}; \quad (18)$$

In other words, to any 3-uplet (Q, K, u) , the computer code associates one square-integrable field from \mathbb{R}^2 to \mathbb{R} . Based on reality constraints and guided by the expert knowledge, the stochastic parameters Q , K , and u of the model are assumed to be all independent with uniform distribution $\mathcal{U}(0, 10)$. For two pollution concentrations C_1 and C_2 with domain in the ground level (in \mathbb{R}^2),

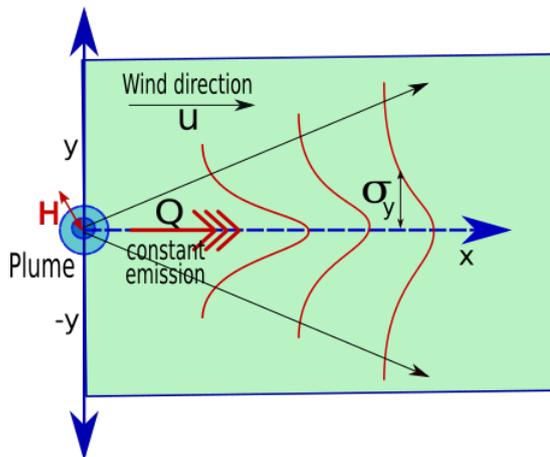


Figure 2: Cross section at $z = 0$ of a contaminant plume emitted from a continuous point source, with wind direction aligned with the x -axis.

the distance used is the classical L^2 distance

$$d(C_1, C_2) = \sqrt{\iint (C_1(x, y, 0) - C_2(x, y, 0))^2 dx dy}.$$

To quantify the sensitivity on the contaminant concentration with respect to Q , K , and u , we consider the family of functions T_a given by $T_{(a_1, a_2)}(b) = \mathbb{1}_{b \in B_{(a_1, a_2)}}$, where a_1 , a_2 , and b square-integrable are applications from \mathbb{R}^2 to \mathbb{R} and $B_{(a_1, a_2)}$ stands for the L^2 ball centered at a_1 with diameter $\overline{a_1 a_2}$ (whence $m = 2$). The values of the indices are presented in Table 1. In this study, we have considered several values of the altitude plume parameter H from 1 to 20 and a sample size N equal to 1000, 2000, and 5000. We observe that, as H increases, the values of the sensitivity indices decrease. When $N = 5000$, we may also observe that the rank of the indices largely varies with respect to the value of H : for large values of H , the parameter K appears to be the most influent on the concentration. In contrast, when $H = 1$, all three parameters seem to have the same influence.

	N=1000			N=2000			N=5000		
	K	Q	u	K	Q	u	K	Q	u
H=1	0.1365	0.1216	0.1330	0.1124	0.1419	0.1453	0.1425	0.1431	0.1562
H=2	0.1028	0.1197	0.1212	0.1291	0.1317	0.1171	0.1222	0.1627	0.1143
H=10	0.0813	0.0891	0.1010	0.1081	0.1077	0.1256	0.0893	0.0831	0.1001
H=20	0.1027	0.0246	0.1041	0.0620	0.0942	0.1030	0.0913	0.0091	0.0329

Table 1: Sensitivity indices for the plume model (17)

4.3 Singular value decomposition in partial differential equation

In this example, we study the sensitivity of the solution (numerical approximation) of an equation in partial derivatives, when the parameters of the equation (inputs) vary. In particular, we analyze the sensitivity of the subspaces generated by the singular value decomposition of the

numerical grid output matrix solution. Following the same example, an elliptical differential partial equation of type diffusive-advective transport is considered:

$$B \frac{\partial C}{\partial t} = \frac{\partial}{\partial x} \left[D \frac{\partial C}{\partial x} \right] + \frac{\partial}{\partial y} \left[D \frac{\partial C}{\partial y} \right] - rC + p_{xy}, \quad (19)$$

with production rate p_{xy} at location (x, y) , consumption rC , and diffusive transport D of a substance C in three dimensions (t, x, y) . The boundaries are prescribed as zero-gradient (default value). The parameter p_{xy} is zero everywhere except for 50 randomly positioned spots denoted by (x_i, y_i) , for $i = 1, \dots, 50$.

Many problems can be modelled by an elliptical differential partial equation. For instance, in physics, electric potential, potential flow, structural mechanics are all studied, see [22]. In biology, the reaction–diffusion–advection equation is used to model chemotaxis observed in bacteria, population migration and evolutionary adaptation to changing environments, see [26].

In this setting, it is usual to compact the information through the singular value decomposition of the solution matrix of dimension 50×50 , that is, the numerical solution of the differential equation (19). Furthermore, it can also be useful to analyze the influence of the parameters in this information compactification. In that view, we assume that the production rate is the same at any of the 50 locations and equal to p and we consider that the function f in (1) defining the output C is given by

$$\begin{aligned} f: \quad \mathbb{R}^4 &\rightarrow L^2(\mathbb{R}_+ \times \mathbb{R}^2) \\ (B, D, r, p) &\mapsto f(B, D, r, p) = (C(t, x, y))_{(t, x, y) \in \mathbb{R}_+ \times \mathbb{R}^2}. \end{aligned} \quad (20)$$

All the input parameters are then assumed to be uniformly distributed:

$$\begin{aligned} B &\sim \mathcal{U}(1 - \beta, 1 + \beta), \\ D &\sim \mathcal{U}(2 - \delta, 2 + \delta), \\ r &\sim \gamma \cdot \mathcal{U}(1, 2), \\ p &\sim \mathcal{U}(0, 1). \end{aligned}$$

Let $C(0, x, y)$ be the solution of (19) and M be the orthogonal matrix given by its singular value decomposition (the first two vectors orthonormal eigenvectors of $C \times C^\top$). Here, the similarity between two matrices is given by the Frobenius distance, that is, for any matrices A_1 and $B_2 \in \mathcal{M}_{n,k}$,

$$d(A_1, A_2) = \sqrt{\text{tr}((A_1 - A_2)^\top (A_1 - A_2))},$$

where $\text{tr}(A)$ represents the trace of the matrix A . The parametric family of functions T_a is given by

$$T_{(a_1, a_2)}(b) = \mathbb{1}_{b \in B_{a_1, a_1 a_2} \cap B_{a_2, a_1 a_2}},$$

where a_1, a_2 and b are matrices and B_{a_1, a_2} still stands for the ball centered at a_1 with diameter $\overline{a_1 a_2}$. In Table 2, the sensitivity indices are calculated for different values of β, δ , and γ and the high influence of the parameter r is observed in all cases. As expected, this influence increases with γ and decreases as the value of δ increases. The simulations have been generated using the R language [18]. In particular, the discretized solution of the differential equation has been computed with the `ReactTran` package [23].

$\gamma = 0.001$	$\delta = 0.1$			$\delta = 0.5$		
	B	D	r	B	D	r
$\beta = 0.1$	0.001	0.011	0.546	0.020	0.071	0.119
$\beta = 0.5$	0.010	0.007	0.491	0.000	0.041	0.102
	$\delta = 0.1$			$\delta = 0.5$		
$\gamma = 0.01$	B	D	r	B	D	r
$\beta = 0.1$	0.000	0.001	0.664	0.010	0.053	0.168
$\beta = 0.5$	0.013	0.006	0.621	0.008	0.041	0.132
	$\delta = 0.1$			$\delta = 0.5$		
$\gamma = 0.1$	B	D	r	B	D	r
$\beta = 0.1$	0.005	0.005	0.794	0.020	0.051	0.179
$\beta = 0.5$	0.000	0.006	0.721	0.000	0.043	0.171

Table 2: Sensitivity indices for the partial differential equation (19)

5 Conclusion

In this paper, we explain how to construct a large variety of sensibility indices as soon as the output space of the black-box model is a general metric space. This construction encompasses the classical Sobol indices [12] and their vectorial generalization [8] as well as some indices based on the whole distribution, namely the Cramér-von-Mises indices [10]. In addition, we propose an estimation procedure that ensures strong consistency and asymptotic normality at a cost of $2N$ calls to the code with a rate of convergence \sqrt{N} . Hence, as soon as $m \geq 1$, this new methodology appears to be more efficient than the so-called Pick-Freeze estimation procedure.

Acknowledgment. We warmly thank Anthony Nouy and Bertrand Iooss for the numerical examples of Sections 4.2 and 4.3. Support from the ANR-3IA Artificial and Natural Intelligence Toulouse Institute is gratefully acknowledged.

References

- [1] E. Borgonovo. A new uncertainty importance measure. *Reliability Engineering & System Safety*, 92(6):771–784, 2007.
- [2] E. Borgonovo, W. Castaings, and S. Tarantola. Moment independent importance measures: New results and analytical test cases. *Risk Analysis*, 31(3):404–428, 2011.
- [3] M. Carrascal, M. Puigcerver, and P. Puig. Sensitivity of gaussian plume model to dispersion specifications. *Theoretical and Applied Climatology*, 48(2-3):147–157, 1993.
- [4] S. Da Veiga. Global sensitivity analysis with dependence measures. *J. Stat. Comput. Simul.*, 85(7):1283–1305, 2015.
- [5] E. De Rocquigny, N. Devictor, and S. Tarantola. *Uncertainty in industrial practice*. Wiley Chisterter England, 2008.
- [6] J.-C. Fort, T. Klein, and N. Rachdi. New sensitivity analysis subordinated to a contrast. *Comm. Statist. Theory Methods*, 45(15):4349–4364, 2016.
- [7] R. Fraiman, F. Gamboa, and L. Moreno. Sensitivity indices for output on a Riemannian manifold. *Accepted in International Journal of uncertainty quantification*, 2019.

- [8] F. Gamboa, A. Janon, T. Klein, and A. Lagnoux. Sensitivity analysis for multidimensional and functional outputs. *Electronic Journal of Statistics*, 8:575–603, 2014.
- [9] F. Gamboa, A. Janon, T. Klein, A. Lagnoux, and C. Prieur. Statistical inference for Sobol pick-freeze Monte Carlo method. *Statistics*, 50(4):881–902, 2016.
- [10] F. Gamboa, T. Klein, and A. Lagnoux. Sensitivity analysis based on Cramér–von Mises distance. *SIAM/ASA J. Uncertain. Quantif.*, 6(2):522–548, 2018.
- [11] W. Hoeffding. A class of statistics with asymptotically normal distribution. *Ann. Math. Statistics*, 19:293–325, 1948.
- [12] A. Janon, T. Klein, A. Lagnoux, M. Nodet, and C. Prieur. Asymptotic normality and efficiency of two Sobol index estimators. *ESAIM: Probability and Statistics*, 18:342–364, 1 2014.
- [13] M. Lamboni, H. Monod, and D. Makowski. Multivariate sensitivity analysis to measure global contribution of input factors in dynamic models. *Reliability Engineering & System Safety*, 96(4):450–459, 2011.
- [14] S. Mahanta, R. Chutia, D. Datta, and H. K. Baruah. Sensitivity analysis with reference to emission concentration of gaussian plume model. *International Journal of Energy, Information and Communications*, 3(2):45–52, 2012.
- [15] A. Owen. Variance components and generalized Sobol’ indices. *SIAM/ASA Journal on Uncertainty Quantification*, 1(1):19–41, 2013.
- [16] A. Owen, J. Dick, and S. Chen. Higher order Sobol’ indices. *Information and Inference*, 3(1):59–81, 2014.
- [17] S. Pouget, M. Bursik, P. Singla, and T. Singh. Sensitivity analysis of a one-dimensional model of a volcanic plume with particle fallout and collapse behavior. *Journal of Volcanology and Geothermal Research*, 326:43–53, 2016.
- [18] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2019.
- [19] A. Saltelli, K. Chan, and E. Scott. *Sensitivity analysis*. Wiley Series in Probability and Statistics. John Wiley & Sons, Ltd., Chichester, 2000.
- [20] I. Sobol. Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. *Mathematics and Computers in Simulation*, 55(1-3):271–280, 2001.
- [21] I. M. Sobol. Sensitivity estimates for nonlinear mathematical models. *Math. Modeling Comput. Experiment*, 1(4):407–414, 1993.
- [22] S. L. Sobolev. *Partial Differential Equations of Mathematical Physics: International Series of Monographs in Pure and Applied Mathematics*. Elsevier, 2016.
- [23] K. Soetaert and F. Meysman. R-package reactran: Reactive transport modelling in r. *Environ. Model. Softw.*, 32:49–60, 2012.
- [24] J. M. Stockie. The mathematics of atmospheric dispersion modeling. *Siam Review*, 53(2):349–372, 2011.

- [25] A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998.
- [26] V. A. Volpert. *Elliptic partial differential equations*, volume 1. Springer, 2011.