



HAL
open science

Reliable error estimation for Sobol' indices

Laurent Gilquin, Lluís Antoni Jiménez Rugama

► **To cite this version:**

Laurent Gilquin, Lluís Antoni Jiménez Rugama. Reliable error estimation for Sobol' indices. 2016.
hal-01358067v1

HAL Id: hal-01358067

<https://inria.hal.science/hal-01358067v1>

Preprint submitted on 30 Aug 2016 (v1), last revised 5 Jan 2017 (v2)

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.

Reliable error estimation for Sobol' indices

Laurent Gilquin · Lluís Antoni Jiménez Rugama

Received: date / Accepted: date

Abstract Insert your abstract here. Include keywords, PACS and mathematical subject classification numbers as needed.

Keywords First keyword · Second keyword · More

1 Introduction

Let f represent a deterministic numerical model in $[0, 1]^d$, $d \geq 1$. Sensitivity measures, also known as Sobol' indices, are used to assess which inputs of f are influential for the output. The normalized indices are scalars between 0 and 1 whose values are interpreted as follows: the closer to 1 the more influential the index. Alternatively, they can be interpreted as the percentage of the variance explained by the inputs. Among all Sobol' indices one can distinguish *first-order* and *total effect* indices. The first measures the effect of a single input, while the latter measures this same effect and includes all its interactions with other inputs.

When dealing with complex numerical models, analytical expressions of Sobol' indices are often inaccessible. In such cases, one must rely on an estimation of these indices. The original estimation procedure is due to Sobol' [13]. However, this procedure requires several model evaluations which are usually expensive. Later on, Saltelli [11] proposed strategies to estimate sets of Sobol' indices at once through the use of a combinatorial formalisms. While elegant, these strategies still require a large number of model evaluations. A cost

efficient alternative to estimate first-order indices was introduced in [7]. This alternative, called replication procedure, has recently been further studied in [14] and generalized to the estimation of closed second-order indices.

A practical problem concerning the use of these methods is how to quantify the number of model evaluations to ensure that Sobol' estimates are accurate enough. This paper addresses this challenge by proposing a reliable error bound for Sobol' indices based on digital sequences. The error bound is defined in terms of the discrete Walsh coefficients of the integrands involved in the Sobol' indices definition. We propose a sequential estimation procedure of Sobol' indices using the error bound as stopping criterion. The procedure operates under the assumption that all integrands lie inside a particular cone of functions (see [3]).

Firstly, section 2 introduces Sobol' indices and reviews both Saltelli's strategy to estimate first-order and total effect Sobol' indices, and the replication procedure. Our main contribution is detailed in Section 3. There we review the construction of the error bound for the estimation of integrals based on digital nets, and then generalize it for Sobol' indices. Section 4 is devoted to analyze the cost in terms of model evaluations of our sequential estimation algorithm. It combines the error bound in Section 3 and either one of the two estimation procedures of Section 2. We also discuss a potential improvement to estimate small first-order indices according to [10]. Finally, examples and illustrations of our procedure are presented in Section 5.

Laurent Gilquin
Inria Grenoble - Rhône-Alpes, Inovallée, 655 avenue de l'Europe,
38330 Montbonnot
E-mail: gilquin.laurent@inria.fr

Lluís Antoni Jiménez Rugama
Illinois Institute of Technology, Rettaliata Engineering Center, 10 W
32st, Chicago, IL 60616
E-mail: ljimene1@hawk.iit.edu

2 Backgrounds on Sobol' indices

2.1 Definition of Sobol' indices

Denote by $\mathbf{x} = (x_1, \dots, x_d)$ the vector of inputs of f and $\mathcal{D} = \{1, \dots, d\}$ the set of dimensions indexes. Let u be a sub-

set of \mathcal{D} , $-u$ its complement and $|u|$ its cardinality. \mathbf{x}_u represents a point in $[0, 1]^{|u|}$ with components $x_j, j \in u$. Given two points \mathbf{x} and \mathbf{x}' , the hybrid point $\mathbf{w} = (\mathbf{x}_u : \mathbf{x}'_{-u})$ is defined as $w_j = x_j$ if $j \in u$ and $w_j = x'_j$ if $j \notin u$. We assume that $f \in \mathbb{L}^2([0, 1]^d)$.

The uncertainty on \mathbf{x} is modeled by a uniform random vector, namely $\mathbf{x} \stackrel{\text{iid}}{\sim} \mathcal{U}([0, 1]^d)$. The Hoeffding decomposition [4, 13] of f is:

$$f(\mathbf{x}) = f_\emptyset + \sum_{u \subseteq \mathcal{D}, u \neq \emptyset} f_u(\mathbf{x}), \quad (1)$$

where:

$$f_\emptyset = \mathbb{E}[f(\mathbf{x})] = \mu, \\ f_u(\mathbf{x}) = \int_{[0, 1]^{|u|}} f(\mathbf{x}) d\mathbf{x}_{-u} - \sum_{v \subsetneq u} f_v(\mathbf{x}).$$

Due to orthogonality, applying the variance to equation (1) leads to the variance decomposition of f :

$$\sigma^2 = \text{Var}[f(\mathbf{x})] = \sum_{u \subseteq \mathcal{D}, u \neq \emptyset} \sigma_u^2, \quad \text{with } \sigma_u^2 = \int_{[0, 1]^{|u|}} f_u(\mathbf{x})^2 d\mathbf{x}_u.$$

From this decomposition, one can express the following two quantities:

$$\underline{\tau}_u^2 = \sum_{v \subseteq u} \sigma_v^2, \quad \bar{\tau}_u^2 = \sum_{v \supsetneq u} \sigma_v^2, \quad u \subsetneq \mathcal{D}.$$

For $u \subsetneq \mathcal{D}$, the two quantities $\underline{\tau}_u^2$ and $\bar{\tau}_u^2$ measure the importance of variables in \mathbf{x}_u . $\underline{\tau}_u^2$ quantifies the main effect of \mathbf{x}_u , that is the effect of all interactions between variables in \mathbf{x}_u . $\bar{\tau}_u^2$ quantifies the main effect of \mathbf{x}_u plus the effect of all interactions between variables in \mathbf{x}_u and variables in \mathbf{x}_{-u} .

$\underline{\tau}_u^2$ and $\bar{\tau}_u^2$ satisfy the following relations: $0 \leq \underline{\tau}_u^2 \leq \bar{\tau}_u^2$ and $\underline{\tau}_u^2 = \sigma^2 - \bar{\tau}_{-u}^2$. These two measures are commonly found in the literature in their normalized form: $\underline{S}_u = \underline{\tau}_u^2 / \sigma^2$ is the closed $|u|$ -order Sobol' index for inputs u , while $\bar{S}_u = \bar{\tau}_u^2 / \sigma^2$ is the total effect Sobol' index of order $|u|$.

The problem of interest is the evaluation of first-order and total effect Sobol' indices. In our framework, we are only interested in single input indices, namely $|u| = 1$. The computation of the normalized indices is performed based on the following integral formulas for their numerators:

$$\underline{\tau}_u^2 = \int_{[0, 1]^{2d-1}} (f(\mathbf{x}) - f(\mathbf{x}_u : \mathbf{x}'_{-u})) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}', \quad (2)$$

$$\bar{\tau}_u^2 = \frac{1}{2} \int_{[0, 1]^{d+1}} (f(\mathbf{x}') - f(\mathbf{x}_u : \mathbf{x}'_{-u}))^2 d\mathbf{x}_u d\mathbf{x}', \quad u \in \mathcal{D}, \quad (3)$$

while variance and mean of f are evaluated as:

$$\sigma^2 = \int_{[0, 1]^d} f(\mathbf{x})^2 d\mathbf{x} - \mu^2, \quad (4) \\ \mu = \int_{[0, 1]^d} f(\mathbf{x}) d\mathbf{x},$$

Usually the complexity of f causes the computation of integrals (2), (3) and (4) to be intractable. In such cases, one can instead estimate these quantities.

2.2 Estimation of Sobol' indices

In this section we review two Monte Carlo procedures for the estimation of Sobol' indices. A design is a point set $\mathcal{P} = \{\mathbf{x}_i\}_{i=0}^{n-1}$ where each point is obtained by sampling each variable x_j n times. Each row of the design is a point \mathbf{x}_i in $[0, 1]^d$ and each column of the design refers to a variable x_j . Consider $\mathcal{P} = \{\mathbf{x}_i\}_{i=0}^{n-1}$ and $\mathcal{P}' = \{\mathbf{x}'_i\}_{i=0}^{n-1}$ two designs where $(\mathbf{x}_i, \mathbf{x}'_i) \stackrel{\text{iid}}{\sim} [0, 1]^{2d}$. One way to estimate the two quantities (2) and (3) is via:

$$\hat{\underline{\tau}}_u^2 = \frac{1}{n} \sum_{i=0}^{n-1} (f(\mathbf{x}_i) - f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u})) f(\mathbf{x}'_i), \quad (5)$$

$$\hat{\bar{\tau}}_u^2 = \frac{1}{2n} \sum_{i=0}^{n-1} (f(\mathbf{x}'_i) - f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u}))^2, \quad u \in \mathcal{D}, \quad (6)$$

using for σ^2 :

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=0}^{n-1} f(\mathbf{x}_i)^2 - \hat{\mu}^2, \quad \text{with } \hat{\mu} = \frac{1}{n} \sum_{i=0}^{n-1} f(\mathbf{x}_i). \quad (7)$$

Then, the Sobol' indices estimators are:

$$\hat{\underline{S}}_u = \hat{\underline{\tau}}_u^2 / \hat{\sigma}^2, \quad \hat{\bar{S}}_u = \hat{\bar{\tau}}_u^2 / \hat{\sigma}^2. \quad (8)$$

The estimation of a single pair $(\hat{\underline{S}}_u, \hat{\bar{S}}_u)$ requires $3n$ evaluations of the model f . Thus, for first order indices, one would need $3nd$ model evaluations. Using a combinatorial formalism, in [11], Saltelli proposes the following estimation strategy:

Theorem 1 *The $d+2$ designs $\{\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u}\}_{i=0}^{n-1}$ constructed for $u \in \{\emptyset, \{1\}, \dots, \{d\}, \mathcal{D}\}$ allows to estimate all first-order and total effect Sobol' indices of order 1 at a cost of $n(d+2)$ evaluations of the model.*

The $d+2$ designs of Theorem 1 are obtained by substituting columns of \mathcal{P} for columns of \mathcal{P}' accordingly to u . Indeed, there is no need of reevaluating $f(\mathbf{x}_i)$ and $f(\mathbf{x}'_i)$ for each u . While elegant, this approach requires a number of model evaluations that grows linearly with respect to the input space dimension.

An efficient alternative to evaluate all first-order indices was proposed by [7] requiring only $2n$ model evaluations. This alternative relies on the construction of two replicated designs. The notion of replicated designs was first introduced by McKay through its replicated Latin Hypercubes in [8]. The definition we provide below generalizes the notion of replicated designs:

Definition 1 Let $\mathcal{P} = \{\mathbf{x}_i\}_{i=0}^{n-1}$ and $\mathcal{P}' = \{\mathbf{x}'_i\}_{i=0}^{n-1}$ be two point sets in $[0, 1]^d$. Let $\mathcal{P}^u = \{\mathbf{x}_{i,u}\}_{i=0}^{n-1}$ (resp. \mathcal{P}'^u), $u \subsetneq \mathcal{D}$, denote the subset of dimensions of \mathcal{P} (resp. \mathcal{P}') indexed by u . We say that \mathcal{P} and \mathcal{P}' are two replicated designs of order $a \in \{1, \dots, d-1\}$ if $\forall u \subsetneq \mathcal{D}$ such that $|u| = a$, \mathcal{P}^u and \mathcal{P}'^u are the same point set in $[0, 1]^a$. We define by π_u the permutation that rearranges the rows of \mathcal{P}'^u into \mathcal{P}^u .

The method introduced in [7] allows to estimate all first-order Sobol' indices with only two replicated designs of order 1. The key point of this method is to use the permutations resulting from the structure of the two replicated designs to mimic the hybrid points in formula (5).

More precisely, let $\mathcal{P} = \{\mathbf{x}_i\}_{i=0}^{n-1}$ and $\mathcal{P}' = \{\mathbf{x}'_i\}_{i=0}^{n-1}$ be two replicated designs of order 1. Denote by $\{f(\mathbf{x}_i)\}_{i=0}^{n-1}$ and $\{f(\mathbf{x}'_i)\}_{i=0}^{n-1}$ the two sets of model evaluations obtained with \mathcal{P} and \mathcal{P}' . From Definition 1, we know that $\mathbf{x}'_{\pi_u(i),u} = \mathbf{x}_{i,u}$. In addition, we define:

$$\pi_u(f(\mathbf{x}'_i)) = f(\mathbf{x}'_{\pi_u(i)}), \quad i \in \{0, \dots, n-1\},$$

then remark:

$$\begin{aligned} \pi_u(f(\mathbf{x}'_i)) &= f(\mathbf{x}'_{\pi_u(i),u} : \mathbf{x}'_{\pi_u(i),-u}), \\ &= f(\mathbf{x}_{i,u} : \mathbf{x}'_{\pi_u(i),-u}). \end{aligned}$$

Hence, each \underline{x}_u^2 can be estimated via formula (5) where $\pi_u(f(\mathbf{x}'_i))$ is exactly $f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u})$ without requiring further model evaluations. This estimation method has been deeply studied and generalized in Tissot et al. [14] to the case of closed second-order indices. In the following we refer to this method as replication procedure.

2.3 Towards a reliable estimation

The aim of this paper is to propose a sequential procedure to estimate first-order and total effect Sobol' indices. A practical problem concerning the estimation of these indices is how large to choose the number of evaluations to ensure that Sobol' estimates are accurate enough. Asymptotical results show that Sobol' estimates are normally distributed ([5, Proposition 2.2], [14, Proposition 3.5]). As a consequence, confidence intervals can be used.

An alternative to address this problematic is the use of sequential procedures. As examples, sequential versions of the replication procedure and McKay's procedure are respectively proposed in [1] and [15]. However, in those two cases, the stopping criterion is a quantity of interest purely empirical, built directly upon the estimates. Such stopping criteria often involve hyper-parameters difficult to tweak but more importantly, fail to guarantee any error bound on the estimates.

Our sequential procedure stands apart from others with the construction of a robust stopping criterion. This criterion is an error bound based on the Walsh series decomposition of the integrands in (2), (3) and (4), and exploits the group properties of digital nets. As such, our procedure relies on an iterative construction of Sobol' sequences. This construction is performed accordingly to the multiplicative approach presented in [2].

The description of the error bound is introduced in the following section and our sequential procedure is detailed in section 4.

3 Reliable error bound for Sobol' indices

We start by reviewing the construction of the error bound proposed in [3] for the estimation of d -dimensional integrals. Then, we present an extension of this error bound for normalized Sobol' indices. This extension is built upon the integral formula of a Sobol' index.

3.1 Reliable integral estimation using digital sequences

We assume we have an embedded sequence of digital nets in base b as in [3, Sec. 2-3],

$$\mathcal{P}_0 = \{\mathbf{0}\} \subset \dots \subset \mathcal{P}_m = \{\mathbf{x}_i\}_{i=0}^{b^m-1} \subset \dots \subset \mathcal{P}_\infty = \{\mathbf{x}_i\}_{i=0}^\infty.$$

Each \mathcal{P}_m has a group structure under the digit wise addition:

$$\mathbf{x} \oplus \mathbf{t} = \left(\sum_{\ell=1}^{\infty} [(x_{j\ell} + t_{j\ell}) \bmod b] b^{-\ell} \pmod{1} \right)_{j=1}^d,$$

where $x_{j\ell}$ and $t_{j\ell}$ are the b -adic decompositions of the j^{th} component of points \mathbf{x} and \mathbf{t} .

In order to relate the group structure of \mathcal{P}_m with the integration error, we introduce the *dual net* which establishes the relationship between any digital net and the *wavenumber* space of non-negative integers \mathbb{N}_0^d . Hence, a dual net is

$$\begin{aligned} \mathcal{P}_m^\perp &= \{\mathbf{k} \in \mathbb{N}_0^d : \langle \mathbf{k}, \mathbf{x} \rangle = 0, \mathbf{x} \in \mathcal{P}_m\}, \\ \langle \mathbf{k}, \mathbf{x} \rangle &= \sum_{j=1}^d \sum_{\ell=0}^{\infty} k_{j\ell} x_{j,\ell+1} \pmod{b}. \end{aligned}$$

and inherits the same embedded structure as for the digital nets,

$$\mathcal{P}_0^\perp = \mathbb{N}_0^d \supset \dots \supset \mathcal{P}_\infty^\perp = \{\mathbf{0}\}. \quad (9)$$

As shown in [3, Sec. 3], the group structure of the digital nets guarantees the property below affecting any Walsh basis $\varphi_{\mathbf{k}}(\mathbf{x}) = e^{2\pi\sqrt{-1}\langle \mathbf{k}, \mathbf{x} \rangle/b}$,

$$\frac{1}{b^m} \sum_{\mathbf{x} \in \mathcal{P}_m} \varphi_{\mathbf{k}}(\mathbf{x}) = \begin{cases} 1, & \mathbf{k} \in \mathcal{P}_m^\perp, \\ 0, & \mathbf{k} \notin \mathcal{P}_m^\perp. \end{cases} \quad (10)$$

Therefore, considering the Walsh decomposition for any $f \in \mathbb{L}^2([0, 1]^d)$:

$$f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{N}_0^d} \hat{f}_{\mathbf{k}} \varphi_{\mathbf{k}}(\mathbf{x}),$$

and,

$$I(f) = \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x},$$

property (10) leads to

$$\begin{aligned} \left| I(f) - \frac{1}{b^m} \sum_{\mathbf{x} \in \mathcal{P}_m} f(\mathbf{x}) \right| &= \left| \hat{f}_0 - \frac{1}{b^m} \sum_{\mathbf{x} \in \mathcal{P}_m} \sum_{\mathbf{k} \in \mathbb{N}_0^d} \hat{f}_{\mathbf{k}} \phi_{\mathbf{k}}(\mathbf{x}) \right|, \\ &= \left| \sum_{\mathbf{k} \in \mathcal{P}_m^+ \setminus \{\mathbf{0}\}} \hat{f}_{\mathbf{k}} \right| \\ &\leq \sum_{\mathbf{k} \in \mathcal{P}_m^+ \setminus \{\mathbf{0}\}} |\hat{f}_{\mathbf{k}}|. \end{aligned} \quad (11)$$

Based on the size of $|\hat{f}_{\mathbf{k}}|$ and the structure of the dual nets (9), in [3, Sec. 4.1] we proposed an ordering of the wavenumbers $\mathbf{k}(\cdot) : \mathbb{N}_0 \rightarrow \mathbb{N}_0^d$. This ordering will be helpful to address the summation in (11). To simplify notation, $\hat{f}_{\mathbf{k}} = \hat{f}_{\mathbf{k}(\kappa)}$. This mapping leads to the following error bound,

$$\left| I(f) - \frac{1}{b^m} \sum_{\mathbf{x} \in \mathcal{P}_m} f(\mathbf{x}) \right| \leq \sum_{\lambda=1}^{\infty} |\hat{f}_{\lambda b^m}|. \quad (12)$$

However, because the knowledge of the Walsh coefficients $\hat{f}_{\mathbf{k}}$ is not assumed, we will estimate them using the fast transform obtained with the precomputed function values, and refer to them as $\tilde{f}_{m,\kappa}$. Note that for function values evaluated at \mathcal{P}_m , one only generates b^m discrete Walsh coefficients $\tilde{f}_{m,\kappa}$.

For $\ell, m \in \mathbb{N}_0$ and $\ell \leq m$ we introduce the following notation,

$$S_m(f) = \sum_{\kappa=[b^{m-1}]}^{b^m-1} |\hat{f}_{\kappa}|, \quad \hat{S}_{\ell,m}(f) = \sum_{\kappa=[b^{\ell-1}]}^{b^{\ell}-1} \sum_{\lambda=1}^{\infty} |\hat{f}_{\kappa+\lambda b^m}|,$$

$$\check{S}_m(f) = \hat{S}_{0,m}(f) + \dots + \hat{S}_{m,m}(f) = \sum_{\kappa=b^m}^{\infty} |\hat{f}_{\kappa}|,$$

$$\tilde{S}_{\ell,m}(f) = \sum_{\kappa=[b^{\ell-1}]}^{b^{\ell}-1} |\tilde{f}_{m,\kappa}|.$$

Finally, we define the set of functions \mathcal{C} ,

$$\mathcal{C} := \{f \in \mathbb{L}^2([0,1]^d) : \hat{S}_{\ell,m}(f) \leq \hat{\omega}(m-\ell) \check{S}_m(f), \ell \leq m, \\ \tilde{S}_m(f) \leq \hat{\omega}(m-\ell) S_{\ell}(f), \ell_* \leq \ell \leq m\}. \quad (13)$$

for $\ell_* \in \mathbb{N}$, $\hat{\omega}$ and $\hat{\omega}$ two non-negative valued functions with $\lim_{m \rightarrow \infty} \hat{\omega}(m) = 0$ and such that $\hat{\omega}(r) \hat{\omega}(r) < 1$ for some $r \in \mathbb{N}$.

Hence, in [3, Sec. 4.2] it is proved that for any $f \in \mathcal{C}$,

$$\left| I(f) - \frac{1}{b^m} \sum_{\mathbf{x} \in \mathcal{P}_m} f(\mathbf{x}) \right| \leq \frac{\tilde{S}_{\ell,m}(f) \hat{\omega}(m) \hat{\omega}(m-\ell)}{1 - \hat{\omega}(m-\ell) \hat{\omega}(m-\ell)} = \varepsilon_{I(f)}, \quad (14)$$

where one may increase m until the error bound $\varepsilon_{I(f)}$ is small enough.

Details concerning the algorithm, the mapping of the wavenumber space, or the meaning and properties of \mathcal{C} , are provided in [3].

For our problem, only Sobol' sequences [12] have been considered. These are digital sequences defined in base $b = 2$. Their major interest comes from their fast and easy implementation as well as their relatively slow growing size rate. Further details concerning Sobol' sequences can be found in [6, 9].

3.2 Extension to Sobol' indices

The idea here is to extend the definition of the error bound (14) to Sobol' indices. To do so, we consider the two integral formulas of the first-order and total effect Sobol' indices:

$$\underline{S}_u(\mathbf{I}) = \frac{\int_{[0,1]^{2d-1}} (f(\mathbf{x}) - f(\mathbf{x}_u : \mathbf{x}'_{-u})) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}'_{-u}}{\int_{[0,1]^d} f(\mathbf{x})^2 d\mathbf{x} - \left(\int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} \right)^2},$$

$$\underline{S}_u(\mathbf{I}) = \frac{I_1}{I_3 - (I_4)^2},$$

$$\bar{S}_u(\mathbf{I}) = \frac{\frac{1}{2} \int_{[0,1]^{d+1}} (f(\mathbf{x}') - f(\mathbf{x}_u : \mathbf{x}'_{-u}))^2 d\mathbf{x}_u d\mathbf{x}'}{\int_{[0,1]^d} f(\mathbf{x})^2 d\mathbf{x} - \left(\int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} \right)^2},$$

$$\bar{S}_u(\mathbf{I}) = \frac{I_2}{I_3 - (I_4)^2},$$

where $\mathbf{I} = (I_1, I_2, I_3, I_4)$ is a vector of integral values. $\underline{S}_u(\mathbf{I})$ and $\bar{S}_u(\mathbf{I})$ are defined as functions over vector \mathbf{I} . If we estimate \mathbf{I} by $\hat{\mathbf{I}}$ with vector of error bounds $\boldsymbol{\varepsilon}_{\mathbf{I}} = (\varepsilon_{I_1}, \varepsilon_{I_2}, \varepsilon_{I_3}, \varepsilon_{I_4})$ according to 3.1, we know that $\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}}) = [\hat{\mathbf{I}} - \boldsymbol{\varepsilon}_{\mathbf{I}}, \hat{\mathbf{I}} + \boldsymbol{\varepsilon}_{\mathbf{I}}]$.

Then, alternatively to the common Sobol' indices estimators (8), we can define the following two estimators with their respective error bounds:

$$\begin{aligned} \hat{\underline{S}}_u &= \frac{1}{2} \left(\min \left(\max_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \underline{S}_u(\mathbf{I}), 1 \right) + \max \left(\min_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \underline{S}_u(\mathbf{I}), 0 \right) \right) \\ \varepsilon_{\hat{\underline{S}}_u} &= \frac{1}{2} \left(\min \left(\max_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \underline{S}_u(\mathbf{I}), 1 \right) - \max \left(\min_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \underline{S}_u(\mathbf{I}), 0 \right) \right) \end{aligned} \quad (15)$$

and,

$$\begin{aligned} \hat{\bar{S}}_u &= \frac{1}{2} \left(\min \left(\max_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \bar{S}_u(\mathbf{I}), 1 \right) + \max \left(\min_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \bar{S}_u(\mathbf{I}), 0 \right) \right) \\ \varepsilon_{\hat{\bar{S}}_u} &= \frac{1}{2} \left(\min \left(\max_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \bar{S}_u(\mathbf{I}), 1 \right) - \max \left(\min_{\mathbf{I} \in B_{\boldsymbol{\varepsilon}_{\mathbf{I}}}(\hat{\mathbf{I}})} \bar{S}_u(\mathbf{I}), 0 \right) \right) \end{aligned} \quad (16)$$

Because numerator and denominator are known to be positive, maximizing $\underline{S}_u(\mathbf{I})$ (resp. $\bar{S}_u(\mathbf{I})$) is done through maximizing the numerator I_1 (resp. I_2) and minimizing the denominator $I_3 - I_4^2$. Analogously, to minimize $\underline{S}_u(\mathbf{I})$ (resp.

$\bar{S}_u(\mathbf{I})$ one minimizes the numerator I_1 (resp. I_2) and maximizes the denominator $I_3 - I_4$.

Under the assumption that each integrand of \mathbf{I} is in \mathcal{C} , these new estimators satisfy:

$$\underline{S}_u \in \left[\hat{\underline{S}}_u - \varepsilon_{\underline{S}_u}, \hat{\underline{S}}_u + \varepsilon_{\underline{S}_u} \right], \quad \bar{S}_u \in \left[\hat{\bar{S}}_u - \varepsilon_{\bar{S}_u}, \hat{\bar{S}}_u + \varepsilon_{\bar{S}_u} \right].$$

For the rest of the paper, we will consider $\hat{\underline{S}}_u$ and $\hat{\bar{S}}_u$ as defined in formulas (15) and (16).

4 Sequential estimation procedure

The sequential estimation procedure we propose combines the error bounds $\varepsilon_{\underline{S}_u}$ and $\varepsilon_{\bar{S}_u}$ presented in the previous section with either one of the two estimation strategies of section 2.2: Saltelli's strategy and the replication procedure.

We start by detailing our procedure in the form of an algorithm. Then, we discuss a possible improvement by considering a new estimator recently introduced in [10] for the estimation of first-order indices.

4.1 Sequential algorithm and its cost

Algorithm 1 summarizes the main steps of our sequential procedure. First, one must fix the tolerance $\varepsilon > 0$ for the Sobol' estimates. Then, we set $m = m_0$ and construct the two designs $\mathcal{P}_m = \{\mathbf{x}_i\}_{i=0}^{2^m-1}$ and $\mathcal{P}'_m = \{\mathbf{x}'_i\}_{i=0}^{2^m-1}$ accordingly to the multiplicative approach detailed in [2]. Step 5 corresponds to the recursive scheme of this approach. The choice of m_0 must be large enough to ensure that each integrand of \mathbf{I} (Section 3.2) is in \mathcal{C} (see (13)).

With this construction, \mathcal{P}_m and \mathcal{P}'_m correspond to the first 2^m points of two independent Sobol' sequences. Furthermore, they possess the structure of two replicated designs of order 1 (Definition 1).

\mathcal{P}_m and \mathcal{P}'_m can be used with Saltelli's strategy to estimate all first-order indices and total effect Sobol' indices. This option is referred as Variant A in Algorithm 1.

Alternatively, \mathcal{P}_m and \mathcal{P}'_m can be used with the replication procedure to estimate all first-order Sobol' indices. This option is referred as Variant B in Algorithm 1.

In both cases we test if the respective error bounds $\varepsilon_{\underline{S}_u}$ and $\varepsilon_{\bar{S}_u}$ are lower than the tolerance ε . For Variant A the stopping criterion writes: $\forall u \in \mathcal{D}, \varepsilon_{\underline{S}_u} \leq \varepsilon$ and $\varepsilon_{\bar{S}_u} \leq \varepsilon$, for Variant B it writes: $\forall u \in \mathcal{D}, \varepsilon_{\underline{S}_u} \leq \varepsilon$.

If the stopping criterion is satisfied, the algorithm stops and Sobol' estimates are returned. Otherwise, m is incremented by one to perform a new estimation.

The cost of our algorithm varies whether Variant A or Variant B is selected. To discuss this cost we note by m^*

Algorithm 1 Sequential estimation of Sobol' indices

```

1: choose  $\varepsilon > 0$ 
2: set:  $m \leftarrow m_0$ 
3:  $bool \leftarrow false$ 
4: while  $!bool$  do
5:    $\mathcal{P}_m \leftarrow \mathcal{P}_{m-1} \cup B_m$ 
    $\mathcal{P}'_m \leftarrow \mathcal{P}'_{m-1} \cup B'_m$ 
6:   for  $u = 1, \dots, d$  do
7:     if Variant A then
8:       estimate  $\hat{\underline{S}}_u$  and  $\hat{\bar{S}}_u$  with formulas (15) and (16) and
       Saltelli's strategy
9:        $bool_u \leftarrow \varepsilon_{\underline{S}_u} \leq \varepsilon \ \& \ \varepsilon_{\bar{S}_u} \leq \varepsilon$ 
10:    end if
11:    if Variant B then
12:      estimate  $\hat{\underline{S}}_u$  with formula (15) and the replication pro-
       cedure
13:       $bool_u \leftarrow \varepsilon_{\underline{S}_u} \leq \varepsilon$ 
14:    end if
15:  end for
16:   $bool \leftarrow \forall u : bool_u$ 
17:   $m \leftarrow m + 1$ 
18: end while
19: return the Sobol' estimates.

```

the ending iteration. If Variant A is selected, the cost of our algorithm writes:

$$\sum_{u \in \mathcal{D}} 2^{m_u} + 2 \times 2^{m^*}, \quad m^* = \max_{u \in \mathcal{D}} m_u,$$

where:

- 2^{m_u} is the number of evaluations $f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u})$ used to estimate both the first-order index \underline{S}_u and the total effect index \bar{S}_u ,
- $2 \times 2^{m^*}$ is the number of evaluations $f(\mathbf{x}_i)$ and $f(\mathbf{x}'_i)$ used in the estimation of each first-order and total effect indices.

If all m_u are equal, the cost of Variant A becomes $2^{m^*}(d+2)$ and we recover the cost specified in Theorem 1 with $n = 2^{m^*}$.

If Variant B is selected the cost of our algorithm equals $2 \times 2^{m^*}$. This cost corresponds to the one of the replication procedure introduced in Section 2.2, where $2n = 2 \times 2^{m^*}$ independent of d .

4.2 Improvement

We focus on the use of a new estimator to evaluate small first-order Sobol' indices in Variant A. This estimator has recently been introduced by Owen in [10] and is called "Correlation 2". Owen discussed and highlighted the efficiency of "Correlation 2" when estimating small first-order indices. Our aim is to show that using "Correlation 2" in Variant A may reduce the total number of model evaluations. Its formula writes as follows:

$$\hat{\underline{t}}_u^2 = \frac{1}{n} \sum_{i=0}^{n-1} (f(\mathbf{x}_i) - f(\mathbf{z}_{i,u} : \mathbf{x}_{i,-u})) (f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u}) - f(\mathbf{x}'_i)),$$

where $(\mathbf{x}_i, \mathbf{x}'_i, \mathbf{z}_i) \stackrel{\text{iid}}{\sim} [0, 1]^{3d}$. It uses an extra set of n model evaluations to estimate $\underline{\tau}_u^2$.

We discuss below the potential improvement brought by using ‘‘Correlation 2’’ in Variant A. The idea is to replace the current estimator (5) by ‘‘Correlation 2’’ for each small first-order indices.

Assume that the number of small first-order indices is known and equals γ . We denote by u_1, \dots, u_γ the indices of the corresponding inputs and $\Gamma = \{1, \dots, \gamma\}$. The cost of Variant A including ‘‘Correlation 2’’ writes:

$$\sum_{j \in \Gamma} 2^{m''_{u_j}} + \sum_{j \in \Gamma} 2^{m'_{u_j}} + \sum_{j \in \mathcal{D} \setminus \Gamma} 2^{m_{u_j}} + 2 \times 2^{m^*}, \quad (17)$$

where:

- for $j \in \Gamma$, $2^{m''_{u_j}}$ is the number of evaluations $f(\mathbf{z}_{i,u_j} : \mathbf{x}_{i,-u_j})$ to estimate \underline{S}_{u_j} ,
- for $j \in \Gamma$, $2^{m'_{u_j}}$ is the number of evaluations $f(\mathbf{x}_{i,u_j} : \mathbf{x}'_{i,-u_j})$ to estimate both \underline{S}_{u_j} and \bar{S}_{u_j} ,
- likewise, for $j \in \mathcal{D} \setminus \Gamma$, $2^{m_{u_j}}$ is the number of evaluations $f(\mathbf{x}_{i,u_j} : \mathbf{x}'_{i,-u_j})$ to estimate both \underline{S}_{u_j} and \bar{S}_{u_j} ,
- $2 \times 2^{m^*}$ is the number of evaluations $f(\mathbf{x}_i)$ and $f(\mathbf{x}'_i)$ used in the estimation of each first-order and total effect index.

Recall that the cost of Variant A without ‘‘Correlation 2’’ writes:

$$\sum_{j \in \Gamma} 2^{m_{u_j}} + \sum_{j \in \mathcal{D} \setminus \Gamma} 2^{m_{u_j}} + 2 \times 2^{m^*}. \quad (18)$$

The difference (17) – (18) equals:

$$\sum_{j \in \Gamma} 2^{m_{u_j}} \left(2^{m''_{u_j} - m_{u_j}} + 2^{m'_{u_j} - m_{u_j}} - 1 \right) = \sum_{j \in \Gamma} c_j. \quad (19)$$

Hence, the sign of this difference indicates whether or not using ‘‘Correlation 2’’ brings an improvement to Variant A. We distinguish two cases :

- 1) for $j \in \Gamma$, the total effect index \bar{S}_{u_j} requires as much or more evaluations than the first-order index \underline{S}_{u_j} . Since the total effect estimator is the same, as a consequence we have $m'_{u_j} = m_{u_j}$ and $c_j > 0$.
- 2) for $j \in \Gamma$, the total effect index \bar{S}_{u_j} requires less evaluations than the first-order index \underline{S}_{u_j} . In this case, if both $m''_{u_j} < m_{u_j}$ and $m'_{u_j} < m_{u_j}$ then $c_j \leq 0$.

Overall we expect to observe case 2) more often than case 1). Indeed, the numerator of \underline{S}_u requires to estimate $2d - 1$ dimensional integrals against only $d + 1$ dimensional integrals for the numerator of \bar{S}_u . Hence, it seems reasonable to expect that it will take less points to estimate \bar{S}_u than \underline{S}_u .

Furthermore, in case 2), we expect the two conditions $m''_{u_j} < m_{u_j}$ and $m'_{u_j} < m_{u_j}$ to always hold as ‘‘Correlation 2’’ is shown in [10] to perform better for small first-order

indices. To support this discussion, illustrations of the use of ‘‘Correlation 2’’ are presented in Section 5.

In practice, one does not know which are the small Sobol’ indices (u_1, \dots, u_γ) . To overcome this issue, we propose the following alternative for Variant A. If at the end of the first iteration, the estimate \hat{S}_u , $u \in \mathcal{D}$, is smaller than a threshold (here 0.1), then at the next iteration estimator (5) is switched for ‘‘Correlation 2’’ and a third Sobol’ sequence $\mathcal{P}''_m = \{\mathbf{z}_i\}_{i=0}^{2^m-1}$ is constructed to obtain the corresponding evaluations $f(\mathbf{z}_{i,u} : \mathbf{x}_{i,-u})$.

5 Applications

5.1 Real case model

5.2 Classical test functions

Acknowledgements This work is supported by the CITIES project funded by the Agence Nationale de la Recherche (grant ANR-12-MONU-0020) and by the United States National Science Foundation (grant DMS-1522687).

The authors thank Fred J. Hickernell and Clémentine Prieur for initiating this collaborative work and Elise Arnaud for her proofreading. The authors are grateful to Stephen Joe, Frances Y. Kuo and Art B. Owen for their helpful answers and suggestions.

References

1. Gilquin, L., Arnaud, E., Monod, H., Prieur, C.: Recursive estimation procedure of Sobol’ indices based on replicated designs, preprint available at <https://hal.inria.fr/hal-01291769>, (2016).
2. Gilquin, L., Jiménez Rugama, L.A., Arnaud, E., Hickernell, F.J., Monod, H., Prieur, C.: Iterative construction of replicated designs based on Sobol’ sequences, preprint available at <https://hal.inria.fr/hal-01349444> (2016)
3. Hickernell, F.J., Jiménez Rugama, L.A.: Reliable Adaptive Cubature Using Digital Sequences: Monte Carlo and Quasi-Monte Carlo Methods, vol. 163, 367-383 (2016)
4. Hoeffding, W.F.: A class of statistics with asymptotically normal distribution, *Ann. Math. Stat.* **19**(3), 293-325 (1948)
5. Janon, A., Klein, T., Lagnoux A., Nodet, M., Prieur C.: Asymptotic normality and efficiency of two Sobol’ index estimators, *ESAIM Probab. Stat.* **18**, 342-364 (2014)
6. Lemieux, C.: Monte Carlo and Quasi-Monte Carlo Sampling, Springer, New-York (2009)
7. Mara, T.A., Rakoto Joseph, O.: Comparison of some efficient methods to evaluate the main effect of computer model factors, *J. Statist. Comput. Simulation* **78**(2), 167-178 (2008)
8. McKay, M.D.: Evaluating prediction uncertainty, Los Alamos National Laboratory Report NUREG/CR- 6311, LA-12915-MS. (1995)
9. Niederreiter, H.: Random Number Generation and Quasi-Monte Carlo Methods: CBMS-NSF Regional Conference Series in Applied Math., vol. 63, SIAM, Philadelphia (1992)
10. Owen, A.B.: Better estimation of small Sobol’ sensitivity indices, *ACM Trans. Model. Comput. Simul.* **23**(2), :11 (2013)

11. Saltelli, A.: Making best use of models evaluations to compute sensitivity indices, *Comput. Phys. Commun.* **145**(2), 280-297 (2002)
12. Sobol', I.M.: On the distribution of points in a cube and the approximate evaluation of integrals, *USSR Comput. Math. Math. Phys.* **7**(4), 86-112 (1967)
13. Sobol', I.M.: Sensitivity indices for nonlinear mathematical models, *Mathematical Modeling and Computational Experiment* **1**, 407-414 (1993)
14. Tissot, J.Y., Prieur, C.: A randomized Orthogonal Array-based procedure for the estimation of first- and second-order Sobol' indices, *J. Statist. Comput. Simulation* **85**(7), 1358-1381 (2015)
15. Tong C.: Self-validated variance-based methods for sensitivity analysis of model outputs, *Reliab. Eng. Syst. Saf.* **95**(3), 301-309 (2010)