

Minimizing the Number of Function Evaluations to Estimate Sobol' Indices Using Quasi-Monte Carlo

Lluís Antoni Jiménez Rugama

Joint work with: Fred J. Hickernell (IIT),
Clémentine Prieur (Univ. Grenoble), Elise Arnaud (Univ. Grenoble),
Hervé Monod (INRA), and Laurent Gilquin (Univ. Grenoble)

Room 120, Bldg E1, Department of Applied Mathematics
Illinois Institute of Technology, Chicago, 60616 IL
Email: ljimene1@hawk.iit.edu

Thursday 26th May, 2016



Outline

- ▶ Introduction
 - ▶ ANOVA
 - ▶ Sobol' Indices
- ▶ Quasi-Monte Carlo Methods
- ▶ Replicated Method



Outline

- ▶ **Introduction**
 - ▶ **ANOVA**—The ANalysis Of VAriance decomposition.
 - ▶ Sobol' Indices
- ▶ Quasi-Monte Carlo Methods
- ▶ Replicated Method



ANOVA

For $f \in L^2([0, 1]^d)$, and $\mathcal{D} = \{1, \dots, d\}$,

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

where,

$$f_u(\mathbf{x}) = \int_{[0,1]^{d-|u|}} f(\mathbf{x}) d\mathbf{x}_{-u} - \sum_{v \subset u} f_v(\mathbf{x}).$$

- ▶ $|u|$ the cardinality of u .
- ▶ $-u := u^c = \mathcal{D} \setminus u$.



Variance Decomposition

Under the previous definitions,

$$\sigma_{\emptyset}^2 = 0, \quad \sigma_u^2 = \int_{[0,1]^d} f_u(\mathbf{x})^2 d\mathbf{x}, \quad \sigma^2 = \int_{[0,1]^d} (f(\mathbf{x}) - \mu)^2 d\mathbf{x}.$$

The ANOVA identity is,

$$\sigma^2 = \sum_{u \subseteq \mathcal{D}} \sigma_u^2.$$



Outline

- ▶ **Introduction**
 - ▶ ANOVA
 - ▶ **Sobol' Indices**—Measuring the importance of each input.
- ▶ Quasi-Monte Carlo Methods
- ▶ Replicated Method



Sobol' Indices

Sobol' introduced the *global sensitivity* indices which measure the variance explained by any dimension subset $u \in \mathcal{D}$:

$$\underline{\tau}_u^2 = \sum_{\substack{v \subseteq u \\ v \in \mathcal{D}}} \sigma_v^2, \quad \text{and} \quad \overline{\tau}_u^2 = \sum_{\substack{v \cap u \neq \emptyset \\ v \in \mathcal{D}}} \sigma_v^2.$$

We have the following properties,

- ▶ $\underline{\tau}_u^2 \leq \overline{\tau}_u^2$.
- ▶ $\underline{\tau}_u^2 + \overline{\tau}_{-u}^2 = \sigma^2$.



Sobol' Indices - Probabilistic Framework

For $\mathbf{X} \sim U[0, 1]^d$, Sobol' indices can also be presented in the following form,

$$\underline{\tau}_u^2 = \text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_u)] = \text{Var} (f(\mathbf{X})) - \mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_u)],$$

$$\overline{\tau}_u^2 = \text{Var} (f(\mathbf{X})) - \text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_{-u})] = \mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_{-u})].$$



The Normalized Sobol' Indices

One may also use the normalized definition of the Sobol' indices,

$$S_u = \frac{\tau_u^2}{\sigma^2} = \frac{\text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_u)]}{\text{Var} (f(\mathbf{X}))} = 1 - \frac{\mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_u)]}{\text{Var} (f(\mathbf{X}))},$$
$$S_u^{\text{tot}} = \frac{\bar{\tau}_u^2}{\sigma^2} = 1 - \frac{\text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_{-u})]}{\text{Var} (f(\mathbf{X}))} = \frac{\mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_{-u})]}{\text{Var} (f(\mathbf{X}))}.$$

satisfying $0 \leq S_u \leq S_u^{\text{tot}} \leq 1$.



The Normalized Sobol' Indices

One may also use the normalized definition of the Sobol' indices,

$$S_u = \frac{\tau_u^2}{\sigma^2} = \frac{\text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_u)]}{\text{Var} (f(\mathbf{X}))} = 1 - \frac{\mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_u)]}{\text{Var} (f(\mathbf{X}))},$$

$$S_u^{\text{tot}} = \frac{\bar{\tau}_u^2}{\sigma^2} = 1 - \frac{\text{Var} [\mathbb{E} (f(\mathbf{X}) | \mathbf{X}_{-u})]}{\text{Var} (f(\mathbf{X}))} = \frac{\mathbb{E} [\text{Var} (f(\mathbf{X}) | \mathbf{X}_{-u})]}{\text{Var} (f(\mathbf{X}))}.$$

satisfying $0 \leq S_u \leq S_u^{\text{tot}} \leq 1$. More specifically, S_u is composed by,

$$S_u = 1 - \frac{I^{(1)}}{I^{(2)} - (I^{(3)})^2}, \quad \text{where} \quad \begin{cases} I^{(1)} \text{ is a } 2d - |u| \text{ dim. integral.} \\ I^{(2)} \text{ is a } d \text{ dim. integral.} \\ I^{(3)} \text{ is a } d \text{ dim. integral.} \end{cases}$$

Error bounds for S_u require more care than error bounds for $I^{(k)}$.



Outline

- ▶ Introduction
 - ▶ ANOVA
 - ▶ Sobol' Indices
- ▶ **Quasi-Monte Carlo Methods**—How can we compute high dimensional integrals efficiently?
- ▶ Replicated Method



Why Quasi-Monte Carlo?

To estimate S_u we need to approximate $I^{(1)}$, $I^{(2)}$, and $I^{(3)}$. However, in high dimensions we need a suitable technique:

Method	Convergence
Trapezoidal rule:	$\mathcal{O}(n^{-2/d})$
Simpson's rule:	$\mathcal{O}(n^{-4/d})$
IID Monte Carlo:	$\mathcal{O}(n^{-1/2})$
Quasi-Monte Carlo:	$\mathcal{O}(n^{-1+\varepsilon})$

(n : number of data points)



Estimating $I^{(1)}$, $I^{(2)}$, and $I^{(3)}$ automatically

Given ε_a and $\mathbf{x} \mapsto f(\mathbf{x})$, we want \hat{I} such that

$$\left| \int_{[0,1]^d} f(\mathbf{x}) \, d\mathbf{x} - \hat{I}(\mathbf{x} \mapsto f(\mathbf{x}), \varepsilon_a) \right| \leq \varepsilon_a,$$

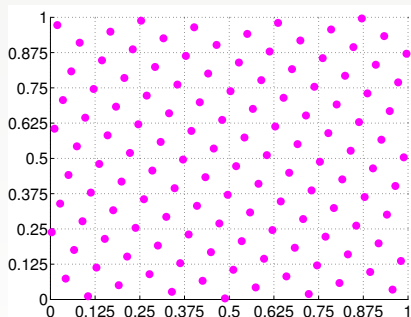
where

$$\hat{I}(\mathbf{x} \mapsto f(\mathbf{x}), \varepsilon_a) = \frac{1}{2^m} \sum_{i=0}^{2^m-1} f(\mathbf{z}_i),$$

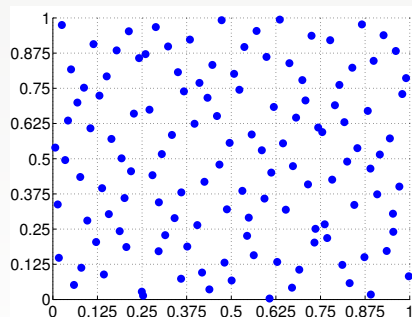
for some **automatic** and **adaptive** choice of m and $\{\mathbf{z}_i\}_{i=0}^{\infty} \in \left\{ \begin{array}{l} \text{Lattice} \\ \text{Digital} \end{array} \right\}$ sequence.



Examples of Sequences



Shifted rank-1 lattice sequence with generating vector (1, 47).



Digitally shifted scrambled Sobol' sequence.



Outline

- ▶ Introduction
 - ▶ ANOVA
 - ▶ Sobol' Indices
- ▶ Quasi-Monte Carlo Methods
- ▶ **Replicated Method**—Reducing the number of function evaluations to compute *first-order* indices.



Normalized *First-Order* Sobol' Indices

In this particular case, we consider $|u| = 1$ and want to estimate $S_u = \sigma_u^2 / \sigma^2$. For this purpose, given $\mathbf{x}, \mathbf{x}' \in [0, 1]^d$, we define the following point,

$$(\mathbf{x}_u : \mathbf{x}'_{-u}) := (x'_1, \dots, x'_{u-1}, x_u, x'_{u+1}, \dots, x'_d) \in [0, 1]^d.$$

Thus, one can use the following integral form to build an estimator:

$$S_u = 1 - \frac{\int_{[0,1]^{2d-1}} \overbrace{f(\mathbf{x})}^{g(\mathbf{x}, \mathbf{x}')} (f(\mathbf{x}) - \overbrace{f(\mathbf{x}_u : \mathbf{x}'_{-u})}^{g_u(\mathbf{x}, \mathbf{x}')}) d\mathbf{x} d\mathbf{x}'_{-u}}{\int_{[0,1]^d} f(\mathbf{x})^2 d\mathbf{x} - \left(\int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} \right)^2} = H(g, g_u).$$



Number of Function Evaluations

We will focus on reducing the number of function evaluations, and to estimate σ_u^2/σ^2 , only g and g_u are evaluated.

Computing all the indices one by one, if one requires n points for each estimation, the total number of function evaluations of g and g_u are

$$2dn,$$

However, if all indices are computed together, g only needs to be evaluated once. Therefore, the number of function evaluations becomes

$$(1 + d)n,$$

Finally, under a special set of quasi-Monte Carlo sequences, this number is decreased to

$$2n.$$



Replicated Designs

Functions g and g_u only share input dimension u :

$$g(\mathbf{x}, \mathbf{x}') = f(x_1, \dots, x_{u-1}, x_u, x_{u+1}, \dots, x_d),$$

$$g_u(\mathbf{x}, \mathbf{x}') = f(x'_1, \dots, x'_{u-1}, x_u, x'_{u+1}, \dots, x'_d).$$

Hence, we can construct our points \mathbf{x}'_i as follows,

$$\begin{pmatrix} x_{0,1} & \cdots & x_{0,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \\ \vdots & & \vdots \end{pmatrix}, \quad \begin{pmatrix} x'_{0,1} & \cdots & x'_{0,d} \\ \vdots & \ddots & \vdots \\ x'_{n,1} & \cdots & x'_{n,d} \\ \vdots & & \vdots \end{pmatrix} = \begin{pmatrix} x_{\pi_1(0),1} & \cdots & x_{\pi_d(0),d} \\ \vdots & \ddots & \vdots \\ x_{\pi_1(n),1} & \cdots & x_{\pi_d(n),d} \\ \vdots & & \vdots \end{pmatrix}.$$



The Right Function Values

Given the right order of points:

$$\begin{pmatrix} \mathbf{x}'_{\pi_u^{-1}(0)} \\ \vdots \\ \mathbf{x}'_{\pi_u^{-1}(n)} \\ \vdots \end{pmatrix} = \begin{pmatrix} x'_{\pi_u^{-1}(0),1} & \cdots & x_{0,u} & \cdots & x'_{\pi_u^{-1}(0),d} \\ \vdots & & \vdots & & \vdots \\ x'_{\pi_u^{-1}(n),1} & \cdots & x_{n,u} & \cdots & x'_{\pi_u^{-1}(n),d} \\ \vdots & & \vdots & & \vdots \end{pmatrix}.$$

Therefore, we only need to evaluate $g_u(\mathbf{x}, \mathbf{x}')$ once:

$$\begin{pmatrix} f(\mathbf{x}'_0) \\ \vdots \\ f(\mathbf{x}'_n) \\ \vdots \end{pmatrix} = \begin{pmatrix} y_0 \\ \vdots \\ y_n \\ \vdots \end{pmatrix} \implies \begin{pmatrix} g_u(\mathbf{x}_0, \mathbf{x}'_0) \\ \vdots \\ g_u(\mathbf{x}_n, \mathbf{x}'_n) \\ \vdots \end{pmatrix} = \begin{pmatrix} y_{\pi_u^{-1}(0)} \\ \vdots \\ y_{\pi_u^{-1}(n)} \\ \vdots \end{pmatrix}$$



Conclusions

- ▶ We can study how each **dimension** explains the overall variance of a model using **Sobol' Indices**.
- ▶ Our **quasi-Monte Carlo automatic cubatures** can be adapted to estimate these indices automatically.
- ▶ **First-order Sobol' Indices** can be estimated using only $2n$ **quasi-Monte Carlo** function evaluations (not depending on d).



References I

Hickernell, F. J. and Li. A. Jiménez Rugama. 2015+. *Reliable adaptive cubature using digital sequences*, Monte Carlo and quasi-Monte Carlo methods 2014. to appear, arXiv:1410.8615 [math.NA].

I.M., Sobol. 2001. *Global sensitivity indices for nonlinear mathematical models and their monte carlo estimates*, Mathematics and Computers in Simulation (MATCOM) **55**, no. 1, 271–280.

Joe, Stephen and Frances Y. Kuo. 2008. *Constructing sobol sequences with better two-dimensional projections*, SIAM J. Scientific Computing **30**, no. 5, 2635–2654.

M. D. McKay, W. J. Conover, R. J. Beckman. 1979. *A comparison of three methods for selecting values of input variables in the analysis of output from a computer code*, Technometrics **21**, no. 2, 239–245.

Niederreiter, H. 1992. *Random number generation and quasi-Monte Carlo methods*, CBMS-NSF Regional Conference Series in Applied Mathematics, SIAM, Philadelphia.

Owen, A. B. 1998. *Scrambling Sobol' and Niederreiter-Xing points*, J. Complexity **14**, 466–489.



References II

Owen, Art B. 2013. *Variance components and generalized Sobol' indices*, SIAM/ASA Journal on Uncertainty Quantification **1**, no. 1, 19–41.

Saltelli, A. 2002. *Making best use of model evaluations to compute sensitivity indices*, Computer Physics Communications **145**, 280–297.

Sobol', I. M. 1967. *The distribution of points in a cube and the approximate evaluation of integrals*, U.S.S.R. Comput. Math. and Math. Phys. **7**, 86–112.

Tissot, Jean-Yves and Clémentine Prieur. 2014. *A randomized Orthogonal Array-based procedure for the estimation of first- and second-order Sobol' indices*, Journal of Statistical Computation and Simulation, 1–24.

