

Automatic estimation of first-order Sobol' indices using the replication procedure

Lluís Antoni Jiménez Rugama, Laurent Gilquin,
Élise Arnaud, Fred J. Hickernell,
Hervé Monod, Clémentine Prieur

Illinois Institute of Technology
Email: ljimene1@hawk.iit.edu

Wednesday 5th July, 2017



Outline

- ▶ **Sobol' Indices**—What are they?
- ▶ Quasi-Monte Carlo Methods
- ▶ Replication Procedure



ANOVA

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

$$f(\mathbf{x}) = \sum_{u \subseteq 1:d} f_u(\mathbf{x}), \quad f_\emptyset = \mu,$$

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 .
- ▶ $-u := u^c = 1:d \setminus u$.

e.g.
$$\underbrace{e^{x_1} \cos(x_2)}_{f(\mathbf{x})} = \underbrace{(e - 1) \sin(1)}_{f_\emptyset} + \underbrace{(e^{x_1} - e + 1) \sin(1)}_{f_{\{1\}}} + \underbrace{(e - 1)(\cos(x) - \sin(1))}_{f_{\{2\}}} + \underbrace{(e^{x_1} - e + 1)(\cos(x) - \sin(1))}_{f_{\{1,2\}}}$$



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,

$$\sigma^2 = \sum_{u \subseteq 1:d} \sigma_u^2.$$



Sobol' indices

Sobol' (1990) and (2001) introduced the *global sensitivity* indices which measure the variance explained by any dimension subset $u \subseteq 1:d$:

$$\underline{\tau}_u^2 = \sum_{\substack{v \subseteq u \\ v \subseteq 1:d}} \sigma_v^2, \quad \text{and} \quad \bar{\tau}_u^2 = \sum_{\substack{v \cap u \neq \emptyset \\ v \subseteq 1:d}} \sigma_v^2.$$

We have the following properties,

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



Normalized closed first-order Sobol' indices

From now on, we consider the normalized Sobol' indices and $|u| = 1$,

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

satisfying $0 \leq S_u \leq 1$. More specifically, S_u is composed by,

$$S_u = 1 - \frac{\mu_1}{\mu_2 - \mu_3^2}, \quad \text{where} \quad \begin{cases} \mu_1 \text{ is a } 2d - 1 \text{ dim. integral.} \\ \mu_2 \text{ is a } d \text{ dim. integral.} \\ \mu_3 \text{ is a } d \text{ dim. integral.} \end{cases}$$

Error bounds for S_u require more care than error bounds for the μ_j .



Improving the estimator

Suppose that $\boldsymbol{\mu} := (\mu_1, \mu_2, \mu_3) \in [\hat{\boldsymbol{\mu}} - \mathbf{err}, \hat{\boldsymbol{\mu}} + \mathbf{err}]$ for

$$\hat{\mu}_j = \frac{1}{n} \sum_{i=0}^{n-1} g_j(\mathbf{x}_i)$$

and some data-based $\hat{\boldsymbol{\mu}}$ and \mathbf{err} . A natural way to estimate S_u is

$$\hat{S}_u = 1 - \frac{\hat{\mu}_1}{\hat{\mu}_2 - \hat{\mu}_3^2},$$

Nonetheless,

$$\tilde{S}_u = 1 - \frac{1}{2} \left(\max_{\boldsymbol{\mu} \in [\hat{\boldsymbol{\mu}} - \mathbf{err}, \hat{\boldsymbol{\mu}} + \mathbf{err}]} \frac{\mu_1}{\mu_2 - \mu_3^2} + \min_{\boldsymbol{\mu} \in [\hat{\boldsymbol{\mu}} - \mathbf{err}, \hat{\boldsymbol{\mu}} + \mathbf{err}]} \frac{\mu_1}{\mu_2 - \mu_3^2} \right),$$

guarantees the tightest absolute error bound.



Why?

For instance, if $\mu \in [1, 1] \times [1, 3] \times [0, 0]$, then

$$\begin{aligned}\frac{1}{3} &\leq S_u \leq 1, \\ |S_u - \hat{S}_u| &= \left| S_u - \frac{1}{2} \right| \leq \frac{1}{2}, \\ |S_u - \tilde{S}_u| &= \left| S_u - \frac{2}{3} \right| \leq \frac{1}{3},\end{aligned}$$

In this case, $1/3$ is the smallest error bound possible.

A deeper study is provided in (Hickernell et al., 2017+) and (Jiménez Rugama and Gilquin, 2017).



Outline

- ▶ Sobol' Indices
- ▶ Quasi-Monte Carlo Methods—How can we estimate S_u efficiently?
- ▶ Replication Procedure



Adaptive quasi-Monte Carlo cubature

To estimate S_u using \tilde{S}_u we need to approximate μ_1 , μ_2 , and μ_3 such that $\mu_j \in [\hat{\mu}_j - \text{err}_j, \hat{\mu}_j + \text{err}_j]$.

Assuming that integrands defining μ_1 , μ_2 , and μ_3 satisfy some specific conditions on the decay of the Fourier coefficients, in (Hickernell and Jiménez Rugama, 2016) and (Jiménez Rugama and Hickernell, 2016) we provided two adaptive quasi-Monte Carlo cubatures that compute $\hat{\mu}_{j,n}$ and $\text{err}_{j,n}$ such that

$$|\mu_j - \hat{\mu}_{j,n}| \leq \text{err}_{j,n} \quad (1)$$

Given the error tolerance ε , the number of points n is increased until

$$\max_{\mu \in [\hat{\mu}_{j,n} - \text{err}_{j,n}, \hat{\mu}_{j,n} + \text{err}_{j,n}]} |S_u(\mu) - \tilde{S}_u| \leq \varepsilon$$



Estimating μ_1 , μ_2 , and μ_3 automatically

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

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

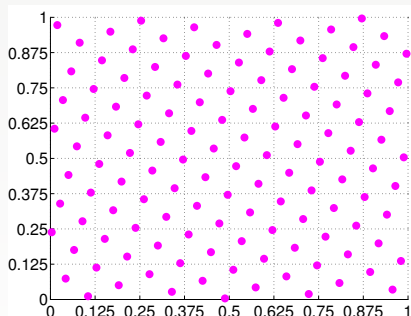
where

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

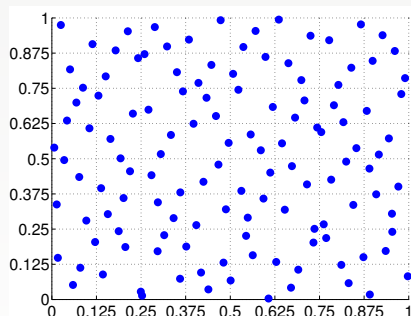
for some **automatic** choice of m and $\{\mathbf{z}_i\}_{i=0}^{\infty} \in \left\{ \begin{array}{l} \text{Lattice} \\ \text{Digital} \end{array} \right\}$ sequence.



Examples of sequences



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

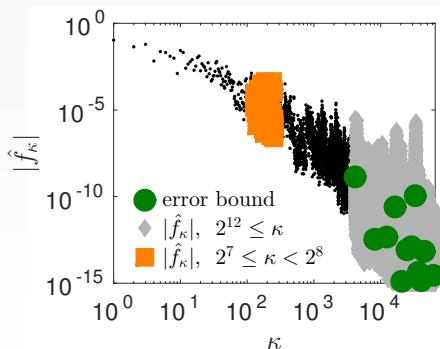


Digitally shifted scrambled Sobol' sequence.



Adaptive algorithm for integrands $f \in \mathcal{C}$

$$\left| \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} - \frac{1}{2^m} \sum_{i=0}^{2^m-1} f(\mathbf{x}_i) \right| \leq \overbrace{\sum \bullet}^{\text{Dual net/lat Fourier coef}} \leq \mathfrak{C}(r, m) \sum_{\kappa=[2^{m-r-1}]}^{2^{m-r}-1} |\tilde{f}_{m,\kappa}| \leq \varepsilon$$



$$\mathcal{C} = \left\{ \begin{array}{l} \sum \text{orange squares} \text{ bounds } \sum \text{grey diamonds} \\ \sum \text{grey diamonds} \text{ bounds } \sum \text{green circles} \end{array} \right\}$$



Outline

- ▶ Sobol' Indices
- ▶ Quasi-Monte Carlo Methods
- ▶ **Replication Procedure**—Reducing the number of function evaluations to compute *first-order* indices.



Number of function evaluations to estimate $\tilde{S}_1, \dots, \tilde{S}_d$

Computing all the indices one by one, if one requires n points for each estimation, the total number of function evaluations is

$$2dn,$$

However, if all indices are computed together, some evaluations can be saved. Therefore, the number of function evaluations becomes

$$(1 + d)n,$$

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

$$2n.$$



Normalized first-order Sobol' indices

Given $\mathbf{x}, \mathbf{x}' \in [0, 1]^d$, we define the following point,

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

This point is used in the definition of S_u :

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



Replicated designs

For the cubature, we must evaluate

$$f(\mathbf{x}_i) = f(x_{i,1}, \dots, x_{i,u-1}, x_{i,u}, x_{i,u+1}, \dots, x_{i,d}),$$

$$f(\mathbf{x}_{i,u}, \mathbf{x}'_{i,-u}) = f(\mathbf{x}'_{i,1}, \dots, \mathbf{x}'_{i,u-1}, x_{i,u}, \mathbf{x}'_{i,u+1}, \dots, \mathbf{x}'_{i,d})$$

for all u . We show that $f(\mathbf{x}_i)$ and $f(\mathbf{x}'_i)$ are enough. We focus on well uniformly distributed points \mathbf{x}_i and \mathbf{x}'_i such that,

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

where the permutations π_u reorder the $\mathbf{x}_{i,u}$ into the $\mathbf{x}'_{i,u}$.



Sobol' points

By construction, Sobol' points have this property when $n = 2^m$. For instance, if $d = 2$ we can use a fourth dimensional Sobol' sequence $\{z_i\}_{i \in \mathbb{N}_0}$:

$$\begin{pmatrix} z_1 & z_2 \\ x_1 & x_2 \\ 0 & 0 \\ 0.5 & 0.5 \\ 0.25 & 0.75 \\ 0.75 & 0.25 \\ 0.125 & 0.625 \\ 0.625 & 0.125 \\ 0.375 & 0.375 \\ 0.875 & 0.875 \\ \vdots & \vdots \end{pmatrix}, \quad \begin{pmatrix} z_3 & z_4 \\ x'_1 & x'_2 \\ 0 & 0 \\ 0.5 & 0.5 \\ 0.25 & 0.75 \\ 0.75 & 0.25 \\ 0.875 & 0.875 \\ 0.375 & 0.375 \\ 0.625 & 0.125 \\ 0.125 & 0.625 \\ \vdots & \vdots \end{pmatrix}, \quad \begin{pmatrix} \pi_1 & \pi_2 \\ 0 & 0 \\ 1 & 1 \\ 2 & 2 \\ 3 & 3 \\ 7 & 7 \\ 6 & 6 \\ 5 & 5 \\ 4 & 4 \\ \vdots & \vdots \end{pmatrix}.$$



No need to evaluate $f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u})$ for d different u

Given the right order of our points \mathbf{x}'_i into \mathbf{x}_i , i.e. π_u^{-1} :

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

Thus, evaluating $f(\mathbf{x}'_i)$, one can directly obtain the $f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u})$:

$$\begin{pmatrix} f(\mathbf{x}'_0) \\ \vdots \\ f(\mathbf{x}'_i) \\ \vdots \end{pmatrix} = \begin{pmatrix} y_0 \\ \vdots \\ y_i \\ \vdots \end{pmatrix} \Rightarrow \begin{pmatrix} f(\mathbf{x}_{0,u} : \mathbf{x}'_{0,-u}) \\ \vdots \\ f(\mathbf{x}_{i,u} : \mathbf{x}'_{i,-u}) \\ \vdots \end{pmatrix} := \begin{pmatrix} y_{\pi_u^{-1}(0)} \\ \vdots \\ y_{\pi_u^{-1}(i)} \\ \vdots \end{pmatrix}$$



Sobol' Indices Example from Bratley et al. (1992)

$$f(\mathbf{x}) = -x_1 + x_1x_2 - x_1x_2x_3 + \cdots + x_1x_2x_3x_4x_5x_6$$

$$\varepsilon = 1\text{E}-3, \quad n = 65\,536$$

u	1	2	3	4	5	6
S_u	0.6529	0.1791	0.0370	0.0133	0.0015	0.0015
\tilde{S}_u	0.6523	0.1796	0.0372	0.0136	0.0015	0.0017
$\hat{S}_u = S_u(\hat{\mu}_n)$	0.6396	0.1787	0.0319	0.0124	0.0000	0.0000



Summary

- ▶ 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.
- ▶ Under some conditions, one may also use **scrambled** Sobol' sequences that keep the **replication property**.
- ▶ The same algorithms can be designed for **rank-1 lattices**.



References I

- Bratley, P., B. L. Fox, and H. Niederreiter. 1992. *Implementation and tests of low-discrepancy sequences* **2**, 195–213.
- Cools, R. and D. Nuyens (eds.) 2016. *Monte Carlo and quasi-Monte Carlo methods: MCQMC, Leuven, Belgium, April 2014*, Springer Proceedings in Mathematics and Statistics, vol. 163, Springer-Verlag, Berlin.
- Gilquin, L., Ll. A. Jiménez Rugama, E. Arnaud, F. J. Hickernell, H. Mond, and C. Prieur. 2016. *Iterative construction of replicated designs based on Sobol' sequences*, C. R. Math. Acad. Sci. Paris **355**, 10–14.
- Hickernell, F. J. and Ll. A. Jiménez Rugama. 2016. *Reliable adaptive cubature using digital sequences*, Monte Carlo and quasi-Monte Carlo methods: MCQMC, Leuven, Belgium, April 2014. arXiv:1410.8615 [math.NA].
- Hickernell, F. J., Ll. A. Jiménez Rugama, and D. Li. 2017+. *Adaptive quasi-monte carlo methods*. Under review.
- Jiménez Rugama, Ll. A. and L. Gilquin. 2017. *Reliable error estimation for Sobol' indices*, Statistics and Computing. in press.



References II

- Jiménez Rugama, Ll. A. and F. J. Hickernell. 2016. *Adaptive multidimensional integration based on rank-1 lattices*, Monte Carlo and quasi-Monte Carlo methods: MCQMC, Leuven, Belgium, April 2014, pp. 407–422. arXiv:1411.1966.
- Owen, Art B. 2013. *Variance components and generalized Sobol' indices*, SIAM/ASA Journal on Uncertainty Quantification **1**, no. 1, 19–41.
- Saltelli, Andrea. 2002. *Making best use of model evaluations to compute sensitivity indices*, Computer Physics Communications **145**, no. 2, 280 –297.
- Sobol', I. M. 1990. *On sensitivity estimation for nonlinear mathematical models*, Matem. Mod. **2**, no. 1, 112–118.
- Sobol', I.M. 2001. *Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates*, Mathematics and Computers in Simulation (MATCOM) **55**, no. 1, 271–280.
- Tissot, J.-Y. and C. Prieur. 2015. *A randomized orthogonal array-based procedure for the estimation of first- and second-order Sobol' indices*, Journal of Statistical Computation and Simulation **85**, no. 7, 1358–1381.

