Institute for Information Transmission Problems of the Russian
Academy of Sciences (Kharkevich Institute)

*as a manuscript*

Panin Ivan Igorevich

# Methods for assessing the accuracy of metamodel-based Sobol' sensitivity indices

PhD Dissertation Summary
for the purpose of obtaining academic degree
Doctor of Philosophy in Computer Science

Moscow – 2023

The PhD dissertation was prepared at *Institute for Information Transmission Problems RAS (Kharkevich Institute)*

Academic Supervisor:    *Doctor of Sciences,*
*Professor,*
*Head of the laboratory at IITP RAS,*
**Vladimir Vyacheslavovich V'yugin**

# Dissertation topic

**Relevance of the research topic.** Sensitivity analysis is an important tool for investigating computational models in engineering and other fields. It allows to find out how various input parameters of the model influence its output and to quantify this effect; in particular, it allows to divide the input parameters into important, relatively important and insignificant [14].

Sensitivity analysis includes a wide range of metrics and techniques including the Morris method, linear regression-based methods, variance-based methods, and others (see review by Iooss [8]). Among all the metrics, we focus on *Sobol' sensitivity indices*, which quantify the portions of the output variance explained by the variability of different input parameters and combinations thereof [16]. The variability of the input parameters of the model is specified using an *a priori* chosen probability distribution. The advantage of the Sobol' method is that it allows to analyze complex nonlinear and nonmonotonic models, and its results are easily interpreted (the effectiveness of the method in engineering is shown in [18]).

The approaches for the evaluation of Sobol' indices are usually divided into *Monte Carlo* and *metamodeling approaches*. Monte Carlo approaches run the analyzed model and conduct high-dimensional numerical integration to estimate Sobol' indices using explicit formulas, such as those given in the articles of Sobol' [15], Owen, and others. The issues of accuracy and convergence of Monte Carlo methods for estimating Sobol' indices are well studied, and it can be concluded that such methods are simple and reliable; however, they require a large number of runs of the analyzed model. On the other hand, *metamodeling approaches* allow one to reduce the required number of model runs. Following these approaches, we replace the original computational model with an approximating *metamodel* (also known as *surrogate model* or *response surface*), which is constructed based on the training set and is better suited for calculating Sobol' indices. For this purpose, metamodels based on Polynomial Chaos Expansion are often used[1], since they allow one to obtain an explicit estimate of Sobol' indices from the expansion coefficients [17].

To date, the issues of accuracy and convergence[2] of metamodel-based Sobol' indices estimates are rather poorly understood, and the existing practical methods for determining their error do not have a rigorous theoretical justification. In particular, in the case of Polynomial Chaos Approximation, the risk bounds for estimates of Sobol' indices are not known, and the only known method for determining the error of corresponding estimates is not justified mathematically [7].

Another range of issues when using the metamodeling approach is related to the choice of the design of experiments (DoE) for building a metamodel. Currently, metamodel-based sensitivity analysis commonly employs DoE techniques that tend to fill the input parameter space "uniformly" in some sense, such as sampling from a dis-

---

[1]The basis of this approximation is polynomials which are orthogonal with respect to the distribution of model input parameters.

[2]With an increase in the size of the training sample used to build the metamodel.

tribution or sampling based on the *latin hypercube*, as well as *quasi-random sequences*. However, such methods do not take into account the features of the metamodel and are not specially adapted for calculating Sobol' indices. It would be expected that the theory of optimal designs would allow the selection of suitable DoE for sensitivity analysis [13]. However, its classical design methods for regression models (*A*-, *D*-, *I*-optimal designs, and others) are associated with maximizing the quality of the approximation itself and do not directly take into account the accuracy of Sobol' indices estimates based on this approximation.

Thus, **it is relevant** to develop both theoretical and practical methods for quality control of estimates of Sobol' indices in the metamodeling approach based on Polynomial Chaos Approximation and effective methods for design of experiments that ensure high accuracy of these estimates.

**The goal of this research** is to create quality control and design of experiments methods for estimating Sobol' sensitivity indices using Polynomial Chaos Approximation. To achieve this goal, the following **research objectives** were identified:

1. Investigate the dependence of the quality of an arbitrary approximation of the analyzed model and the accuracy of estimates of Sobol' indices obtained on the basis of this approximation.

2. Develop an algorithm for estimating the error of Sobol' indices calculated using Polynomial Chaos Approximation.

3. Perform a theoretical analysis of the accuracy of Sobol' indices estimates obtained based on Polynomial Chaos Approximation.

4. Develop a DoE algorithm that is effective for estimating Sobol' indices based on Polynomial Chaos Approximation.

5. Implement the proposed algorithms in the software package.

# Key results

**The main defense points:**

1. A relationship has been established between the error of estimates of Sobol' indices and the error of approximation by which these estimates were obtained. It is shown that the corresponding upper bound on the error of these estimates is achievable.

2. For metamodels of the Polynomial Chaos Approximation type, a method has been developed to control the error of metamodel-based estimates of Sobol' indices, which uses the proven upper bound on the error of these estimates.

3. Non-asymptotic upper bounds on the risk of metamodel-based estimates of Sobol' indices under the condition of random design for Polynomial Chaos Approximation are proved. For these bounds, in the case of specific families of multidimensional orthogonal polynomials, convergence rates are found.

4. An asymptotic distribution of estimates of Sobol' indices is found and a method for sequential DoE is developed, which allows improving the average accuracy of these estimates compared to standard DoE methods.

5. With the help of the developed software package, a number of engineering problems were solved. In particular, its application to solve the problem of analyzing the factors that affect the magnitude of the deflection of a bar structure (truss) under the action of an external load, made it possible to increase the average accuracy of estimates of Sobol' indices by 10% compared to standard DoE methods.

**Scientific novelty.** This study for the first time raises and resolves the question of the relationship between the error of the estimates of Sobol' indices and the theoretical error of an arbitrary approximation, with the help of which these estimates were obtained; which opens up new possibilities for analyzing the accuracy of metamodel-based Sobol' indices estimates using arbitrary types of approximations.

Based on the obtained relationship between the error of indices estimates and the approximation error, an applied method for quality control of Sobol' indices estimates is proposed. This relation also made it possible for the first time to find the upper bounds of the risk of Sobol' indices estimates based on Polynomial Chaos Approximation under the condition of random design.

In addition, the study proposes the idea of using the theory of optimal designs for regression models for more efficient calculation of Sobol' indices. In particular, for the case of Polynomial Chaos Approximation, based on the criterion of $D$-optimality, a new method of sequential design of experiments was developed, which is effective for estimating sensitivity indices. The proposed approach can be extended to other optimality criteria.

**Author's contribution to the study.** The content of the dissertation and the main points, submitted for defense, reflect the personal contribution of the author of the dissertation to the published articles.

The article [10], which analyzes the accuracy of estimates of Sobol' indices, proposes a method for controlling their quality, and obtains risk bounds for them, was written without co-authors.

Preparation for publication of a series of articles [3, 4, 5] devoted to the DoE for sensitivity analysis was carried out jointly with co-authors, and the contribution of the dissertation author was decisive. In this series of articles, E. V. Burnaev proposed a general statement of the problem, general approaches to its solution, and ideas for some experiments. In addition, the proof of Theorem 1 in [5] was obtained jointly with E. V. Burnaev; the rest of the results in this article, including the proposed

DoE algorithm, belong personally to the dissertation author. The program code for finite element computational models for testing the DoE method proposed in [4, 5] was provided by B. Sudret. The development of a software package that implements the proposed methods, and all computational experiments were performed by the author of the dissertation.

**Research methods.** To achieve research objectives, the methods of mathematical statistics, probability theory, approximation theory, matrix algebra and Fourier analysis were used.

**Theoretical and practical significance.** From a theoretical point of view, the results of the dissertation provide a basis for analyzing the accuracy of estimates of Sobol' indices obtained based on various types of approximations and using various designs of experiments. From a practical point of view, the results complement and improve existing approaches in sensitivity analysis of mathematical models in engineering and other fields.

# Publications and approbation of research

### First-tier publications

1. Ivan Panin. "Risk of estimators for Sobol' sensitivity indices based on metamodels". In: *Electron. J. Statist.* 15.1 (2021), pp. 235–281. ISSN: 1935-7524. DOI: 10.1214/20-EJS1793. URL: https://projecteuclid.org/euclid.ejs/1609902190

### Second-tier publications

1. Evgeny Burnaev, Ivan Panin, and Bruno Sudret. "Efficient design of experiments for sensitivity analysis based on polynomial chaos expansions". In: *Ann. Math. Artif. Intell.* 81.1-2 (2017), pp. 187–207. ISSN: 1012-2443. DOI: 10.1007/s10472-017-9542-1. URL: https://doi.org/10.1007/s10472-017-9542-1

2. Evgeny Burnaev, Ivan Panin, and Bruno Sudret. "Effective Design for Sobol Indices Estimation Based on Polynomial Chaos Expansions". In: *Conformal and Probabilistic Prediction with Applications.* Ed. by Alexander Gammerman et al. Cham: Springer International Publishing, 2016, pp. 165–184. ISBN: 978-3-319-33395-3

3. Evgeny Burnaev and Ivan Panin. "Adaptive Design of Experiments for Sobol Indices Estimation Based on Quadratic Metamodel". In: *Statistical Learning and Data Sciences.* Ed. by Alexander Gammerman, Vladimir Vovk, and Harris Papadopoulos. Cham: Springer International Publishing, 2015, pp. 86–95. ISBN: 978-3-319-17091-6

### Other publications

1. Ivan Panin and Pavel Prikhodko. "Approaches to the evaluation of the sensitivity indices variance in the problem of global sensitivity analysis". Russian. In: *Proceedings of the conference "Information Technology and Systems"* (Petrozavodsk, Russia). IITP RAS. 2012, pp. 173–178. ISBN: 978-5-901158-19-7. URL: http://www.itas2012.iitp.ru/pdf/1569602539.pdf

**Reports at conferences and seminars**

1. 5-th International Symposium on Conformal and Probabilistic Prediction with Applications (2016, Madrid, Spain);

2. 3-rd International Symposium on Statistical Learning and Data Sciences (2015, Egham, UK);

3. 37-th Conference for Young Scientists and Engineers "Information Technology and Systems" (2013, Kaliningrad, Russia);

4. 55-th MIPT Scientific Conference (2012, Dolgoprudny, Russia);

5. 35-th Conference for Young Scientists and Engineers "Information Technology and Systems" (2012, Petrozavodsk, Russia);

6. Structural Learning Seminar at IITP RAS (2019, Moscow, Russia).

# Contents

**The Introduction** substantiates the relevance of the dissertation, formulates the research goal and argues the scientific novelty of the study, and also shows the theoretical and practical significance of the dissertation.

**In the first chapter**, we consider the problem statement of global sensitivity analysis, give an example of such a problem from engineering practice, and describe an approach to its solution using Sobol' sensitivity indices. In addition, a method for calculating Sobol' indices using Polynomial Chaos Approximation is considered.

**The first part of the chapter** introduces basic concepts and describes the Sobol' method. Consider a function $y = f(\mathbf{x})$, which corresponds to some physical model. The vector of input variables $\mathbf{x} = (x_1, \ldots, x_d)^T \in \mathscr{X} \subseteq \mathbb{R}^d$ is *model parameters* in the *design space*; and the output variable $y \in \mathbb{R}^1$ is *the model response*. We will assume that this function is a "black box" that can be studied by setting different values of the input parameters and analyzing the response. It is assumed that the calculation of the response can take a long time.

An informal statement of the *global sensitivity analysis* problem is to quantify the "importance" of various model input parameters and their groups, given some kind of model input variability; and, thus, to rank parameters according to the degree of influence on the response of the model, highlighting important and non-essential

parameters. It is also assumed that the input parameters can vary in some wide range, and not in a narrow neighborhood.

One of the ways to formalize the concept of "importance" of parameters is associated with the method of I. M. Sobol'. Assume that on the set $\mathscr{X}$ some *a priori* known probability measure $\mu$ is given, which is a product of measures: $\mu = \otimes_{i=1}^{d}\mu_i$, where $\mu_i$ is a probability measure on $\mathscr{X}_i \subseteq \mathbb{R}$, and the set $\mathscr{X}$ can be represented as Cartesian product $\mathscr{X} = \mathscr{X}_1 \times \ldots \times \mathscr{X}_d$. The corresponding probability distribution describes the uncertainty and/or variability of the input parameters, modeled by a random vector $\mathbf{x} = (x_1, \ldots, x_d)^T$ with independent components. In this setting, the model output $y = f(\mathbf{x})$ also becomes a stochastic variable.

Assume that the function $f$ lies in Hilbert space $L^2(\mathscr{X}, \mu)$ of real-valued functions on $\mathscr{X}$ that are square-integrable with respect to the measure $\mu$. We have the following unique Sobol-Hoeffding decomposition [19] of the model output given by

$$
\begin{aligned}
f(\mathbf{x}) &= f_0 + \sum_{i=1}^{d} f_i(x_i) + \sum_{1 \leq i < j \leq d} f_{ij}(x_i, x_j) + \ldots + f_{1\ldots d}(x_1, \ldots, x_d) \\
&= \sum_{\mathcal{U} \subseteq \{1,\ldots,d\}} f_{\mathcal{U}}(\mathbf{x}_{\mathcal{U}})
\end{aligned}
\tag{1}
$$

with $2^d$ terms of increasing dimension that satisfy the condition

$$
\mathbb{E}_{\mu_i}[f_{\mathcal{U}}] = \int_{\mathscr{X}_i} f_{\mathcal{U}}(\mathbf{x}_{\mathcal{U}}) d\mu_i(x_i) = 0 \quad \text{for} \quad \forall i \in \mathcal{U},
\tag{2}
$$

where $\mathcal{U} \subseteq \{1, 2, \ldots, d\}$ is a subset of indexes of input variables; $\mathbf{x}_{\mathcal{U}}$ — the vector with components $(x_i, \ i \in \mathcal{U})^T$ and $f_{\varnothing} \triangleq f_0 = \mathbb{E}_{\mu}[f(\mathbf{x})]$. The terms of this decomposition are orthogonal:

$$
\mathbb{E}_{\mu}\big[f_{\mathcal{U}}(\mathbf{x}_{\mathcal{U}}) f_{\mathcal{V}}(\mathbf{x}_{\mathcal{V}})\big] = 0, \quad \text{if} \quad \mathcal{U} \neq \mathcal{V}, \quad \text{where} \quad \mathcal{U}, \mathcal{V} \subseteq \{1, 2, \ldots, d\},
$$

from which one can obtain the following decomposition of the variance of model response:

$$
\mathbb{V}_{\mu}[f(\mathbf{x})] = \sum_{\mathcal{U} \subseteq \{1,\ldots,d\}} \mathbb{V}_{\mu}[f_{\mathcal{U}}(\mathbf{x}_{\mathcal{U}})].
\tag{3}
$$

Assuming $\mathbb{V}_{\mu}[f] > 0$, we introduce *Sobol' sensitivity indices* (SI).

**Definition 1.** *The Sobol' index of the set $\mathbf{x}_{\mathcal{U}}$, $\mathcal{U} \subseteq \{1, \ldots, d\}$ of input variables of a function is defined as*

$$
S_{\mathcal{U}} = \frac{\mathbb{V}_{\mu}[f_{\mathcal{U}}(\mathbf{x}_{\mathcal{U}})]}{\mathbb{V}_{\mu}[f(\mathbf{x})]}.
\tag{4}
$$

Denote $\mu_{\mathcal{U}} \triangleq \otimes_{i \in \mathcal{U}} \mu_i$, $\mathbb{E}_{\mathcal{U}} \triangleq \mathbb{E}_{\mu_{\mathcal{U}}}$, $\mathbb{V}_{\mathcal{U}} \triangleq \mathbb{V}_{\mu_{\mathcal{U}}}$, and $\sim\mathcal{U} \triangleq \{1, \ldots, d\} \backslash \mathcal{U}$. Then for

$\mathcal{U} = \{i\}$ the Sobol' index (4) can be represented as

$$S_i = \frac{\mathbb{V}_i\big[\mathbb{E}_{\sim i}[f(\mathbf{x})|x_i]\big]}{\mathbb{V}_\mu[f]}, \quad i = 1, \ldots, d, \tag{5}$$

We also define a quantity that characterizes the "total" contribution of a group of variables to the variability of the model — *total-effect index* (also known as *total Sobol' index* and *total-index*).

**Definition 2.** *The total-effect index of the set* $\mathbf{x}_\mathcal{U}$, $\mathcal{U} \subseteq \{1, \ldots, d\}$ *of input variables of a function is defined as*

$$T_\mathcal{U} = \sum_{\mathcal{U} \cap \mathcal{V} \neq \varnothing} S_\mathcal{V} = \frac{\mathbb{E}_{\sim \mathcal{U}}\big[\mathbb{V}_\mathcal{U}[f(\mathbf{x})|\mathbf{x}_{\sim \mathcal{U}}]\big]}{\mathbb{V}_\mu[f]}. \tag{6}$$

Thus, Sobol' sensitivity indices introduced above offer one of the possible ways to formalize and solve the problem of global sensitivity analysis.

**The second part** is devoted to the method for calculating Sobol' indices and total-effect indices using the metamodeling approach; in particular, with the help of Polynomial Chaos Approximation.

Direct calculation of Sobol' indices leads to computationally expensive multidimensional integration. To simplify this problem using the metamodeling approach, one can replace the original function $f(\mathbf{x})$ with the approximation $\widehat{f}(\mathbf{x})$ that is better suited for computing of Sobol' indices.

We will mainly consider metamodels of the Polynomial Chaos Approximation type, which are often encountered in sensitivity analysis problems. Introduce them formally. Denote scalar product and norm for $g, h \in L^2(\mathscr{X}, \mu)$ as $\langle g, h \rangle_\mu = \int_{\mathbf{x} \in \mathscr{X}} g(\mathbf{x})h(\mathbf{x})d\mu(\mathbf{x})$ and $\|g\|_\mu^2 \triangleq \|g\|_{L^2(\mathscr{X},\mu)}^2 = \int_{\mathbf{x} \in \mathscr{X}} g^2(\mathbf{x})d\mu(\mathbf{x})$; and *supremum norm* as $\|g\|_{L^\infty} \triangleq \|g\|_{L^\infty(\mathscr{X})} \triangleq \sup_{\mathbf{x} \in \mathscr{X}} |g(\mathbf{x})|$. The norm in Euclidean vector spaces is denoted as $\| \cdot \|$.

Assume that there is a function set $\{\Psi_{\boldsymbol{\alpha}}(\mathbf{x})\}$ in $L^2(\mathscr{X}, \mu)$ parameterized by multi-index[3] $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_d) \in \mathbb{N}^d$, which consists of $\mu$-orthonormal functions and has the form of tensor product of $d$ families of $\mu_i$-orthonormal one-dimensional functions $\{\psi_{\alpha_i}^{(i)}, \alpha_i \in \mathbb{N}\}$ with $\psi_0^{(i)} \equiv 1$ and $\mathbb{E}_i\big[\psi_\alpha^{(i)}(x_i)\psi_\beta^{(i)}(x_i)\big] = \delta_{\alpha\beta}$ for $\alpha, \beta \in \mathbb{N}$, where $\delta$ is the Kronecker symbol. As a result,

$$\Psi_{\mathbf{0}}(\mathbf{x}) \equiv 1, \quad \Psi_{\boldsymbol{\alpha}}(\mathbf{x}) = \prod_{i=1}^{d} \psi_{\alpha_i}^{(i)}(x_i), \quad \mathbf{x} \in \mathscr{X},$$

$$\langle \Psi_{\boldsymbol{\alpha}}, \Psi_{\boldsymbol{\beta}} \rangle_\mu = \delta_{\boldsymbol{\alpha}\boldsymbol{\beta}}, \quad \boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{N}^d. \tag{7}$$

Define the metamodel as a linear combination of $N$ functions from the set

---

[3]Denote $\mathbb{N} \triangleq \{0, 1, 2, \ldots\}$, $\mathbb{N}_+ \triangleq \{0, 1, 2, \ldots\}$ and introduce zero multi-index $\mathbf{0} \triangleq (0, \ldots, 0) \in \mathbb{N}^d$.

$\{\Psi_{\boldsymbol{\alpha}}(\mathbf{x}), \ \boldsymbol{\alpha} \in \mathscr{L}_N\}$ for some set of multi-indices $\mathscr{L}_N \subset \mathbb{N}^d$:

$$\widehat{f}(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N} \widehat{c}_{\boldsymbol{\alpha}} \Psi_{\boldsymbol{\alpha}}(\mathbf{x}), \quad \mathbf{x} \in \mathscr{X}, \ \widehat{c}_{\boldsymbol{\alpha}} \in \mathbb{R}. \tag{8}$$

If the regressors $\{\Psi_{\boldsymbol{\alpha}}\}$ are multidimensional polynomials, then (8) corresponds to Polynomial Chaos Approximation[4]. One of the important advantages of the presented type of approximation is that it allows calculating Sobol' indices analytically from the expansion coefficients [12]. Indeed, if $\mathbb{V}_{\mu}[\widehat{f}] > 0$, then the Sobol' index of a non-empty set of input variables $\mathbf{x}_{\mathcal{U}}, \mathcal{U} \subseteq \{1, \ldots, d\}$ of a function $\widehat{f}$ defined as (8) is expressed as

$$\widehat{S}_{\mathcal{U}}(\widehat{\mathbf{c}}) = \frac{\sum_{\boldsymbol{\alpha} \in \mathbb{L}_{\mathcal{U}}} \widehat{c}_{\boldsymbol{\alpha}}^2}{\sum_{\boldsymbol{\alpha} \in \mathscr{L}_N^+} \widehat{c}_{\boldsymbol{\alpha}}^2}, \tag{9}$$

where $\mathscr{L}_N^+ \triangleq \mathscr{L}_N \backslash \mathbf{0}$, and $\mathbb{L}_{\mathcal{U}} \triangleq \mathbb{L}_{\mathcal{U}}[\mathscr{L}_N]$ is the subset of $\mathscr{L}_N$ that consists of such multi-indices that only indices corresponding to variables $\mathbf{x}_{\mathcal{U}}$ are nonzero: $\mathbb{L}_{\mathcal{U}} = \{\boldsymbol{\alpha} \in \mathscr{L}_N^+: \ \alpha_i > 0 \text{ for all } i \in \mathcal{U}; \ \alpha_i = 0 \text{ for } i \notin \mathcal{U}\}$. The vector of coefficients $(\widehat{c}_{\boldsymbol{\alpha}}, \ \boldsymbol{\alpha} \in \mathscr{L}_N)^T$ is denoted as $\widehat{\mathbf{c}} \in \mathbb{R}^N$. Note that $\mathbb{V}_{\mu}[\widehat{f}] = \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N^+} \widehat{c}_{\boldsymbol{\alpha}}^2$. Total-effect indices for $\widehat{f}$ are calculated similarly.

**In the third part**, an observation model is introduced and methods for constructing Polynomial Chaos Approximation for a finite training set are considered.

The observation model. We will assume that the only information about the model $f$ comes from the observations; more precisely, for some *design of experiments* $\mathcal{D} = (\mathbf{x}_i \in \mathscr{X})_{i=1}^n \in \mathbb{R}^{n \times d}$ one can obtain a set of model responses and form *a training sample*:

$$\mathcal{S} = \big(\mathbf{x}_i, \ y_i = f(\mathbf{x}_i) + \eta_i\big)_{i=1}^n, \tag{10}$$

where $\eta_i$ are independent and identically distributed random measurement errors such that $\mathbb{E}\eta_i = 0$, $\mathbb{V}\eta_i = \sigma^2 < \infty$; and are independent from $\mathbf{x}$. In matrix form: $\mathcal{S} = \big(\mathcal{D} \in \mathscr{X}^n, \ Y = f(\mathcal{D}) + \boldsymbol{\eta} \in \mathbb{R}^n\big)$, where $f(\mathcal{D}) \triangleq \big(f(\mathbf{x}_1), \ldots, f(\mathbf{x}_n)\big)^T$ and $\boldsymbol{\eta} = (\eta_1, \ldots, \eta_n)^T$. We will consider both the general case of *noisy* observations and the special *noiseless* case corresponding to $\sigma^2 = 0$.

Building a metamodel. Consider some fixed sequence of nested sets of $d$-dimensional multi-indices:

$$\{\mathbf{0}\} = \mathscr{L}_1 \subset \ldots \subset \mathscr{L}_N \ldots \subset \mathscr{L}_{\infty} = \mathbb{N}^d, \tag{11}$$

where $|\mathscr{L}_N| = N$ for all $N \in \mathbb{N}_+$. Each element of this sequence corresponds to the metamodel of the form (8). We will call each $\mathscr{L}_N$ a *truncation set*, since for the corresponding $N$-th approximation $\widehat{c}_{\boldsymbol{\alpha}} \triangleq 0$ if $\boldsymbol{\alpha} \notin \mathscr{L}_N$. Fix some truncation set $\mathscr{L}_N$, and denote the subspace of all linear combinations of $\{\Psi_{\boldsymbol{\alpha}}, \ \boldsymbol{\alpha} \in \mathscr{L}_N\}$ as $V_N \triangleq span\{\Psi_{\boldsymbol{\alpha}}, \ \boldsymbol{\alpha} \in \mathscr{L}_N\}$. *Theoretical orthogonal projection* of $f$ onto $V_N$ with

---

[4]We can consider not only polynomials. At the same time, this is an important special case.

respect to $\mu$-norm is defined as

$$f_N \triangleq \underset{\widehat{f} \in V_N}{\text{argmin}} \| f - \widehat{f} \|_\mu. \tag{12}$$

We also define $e_N(\mathbf{x}) \triangleq f(\mathbf{x}) - f_N(\mathbf{x})$. The function $f_N$ corresponds to the best possible approximation of the model $f$ in the space $V_N$ with respect to $\mu$-norm.

Consider now the construction of an approximation for a finite training set. Let $\Phi \triangleq \boldsymbol{\Psi}(\mathcal{D}) = \big( \Psi_{\boldsymbol{\alpha}}(\mathcal{D}), \ \boldsymbol{\alpha} \in \mathscr{L}_N \big) \in \mathbb{R}^{n \times N}$ and $\Psi_{\boldsymbol{\alpha}}(\mathcal{D}) \triangleq \big( \Psi_{\boldsymbol{\alpha}}(\mathbf{x}_1), \ldots, \Psi_{\boldsymbol{\alpha}}(\mathbf{x}_n) \big)^T \in \mathbb{R}^n$. Among all methods for estimating the expansion coefficients, we will consider the following:

- Projection method based on quasi-regression [1]

$$\widehat{\mathbf{c}}^P = \frac{1}{n} \Phi^T Y \in \mathbb{R}^N. \tag{13}$$

- Ordinary Least Squares, LS

$$\widehat{\mathbf{c}}^{LS} = (\Phi^T \Phi)^{-1} \Phi^T Y, \text{ if } \det(\Phi^T \Phi) \neq 0. \tag{14}$$

We will refer to the approximations constructed based on these two methods as $\widehat{f}^P$ and $\widehat{f}^{LS}$ correspondingly. Related Sobol' indices, estimated via these two approximations, are denoted by $\widehat{S}^P$ and $\widehat{S}^{LS}$, respectively.

**In the second chapter**, we perform a theoretical analysis of the error of Sobol' indices estimates based on arbitrary metamodels and establish the relationship of this error with the accuracy of metamodels. Based on this analysis, a new quality control method for such estimates is proposed. It is also shown that obtained upper bounds for the error of indices estimates are achievable.

All further results are valid for both Sobol' indices and total-effects of all orders unless otherwise stated. In order to avoid duplication, we use the notation $S_{\mathcal{U}}$ for indices of both types in theorems' statements.

First, we need to make sure that the closeness of the function and its arbitrary[5] approximation $f \approx \widehat{f}$ leads to the closeness of their (total) Sobol' indices $S_{\mathcal{U}}$ and $\widehat{S}_{\mathcal{U}}$. Note that the opposite is not true in general. To characterize the closeness of functions, we will use the relative error of approximation given by

$$\mathscr{E} \triangleq \frac{\| f - \widehat{f} \|_\mu}{\mathbb{V}_\mu^{1/2}[f]}. \tag{15}$$

**Theorem 1.** *For any functions $f, \widehat{f} \in L^2(\mathscr{X}, \mu)$ such that $\mathbb{V}_\mu[f] > 0$, $\mathbb{V}_\mu[\widehat{f}] > 0$, it*

---

[5]In particular, the approximation may not be related to polynomial chaos.

holds for corresponding Sobol' and total-effect indices for $\mathcal{U} \subseteq \{1, \ldots, d\}$

$$\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq \left\{\sqrt{S_{\mathcal{U}}(1 - \widehat{S}_{\mathcal{U}})} + \sqrt{\widehat{S}_{\mathcal{U}}(1 - S_{\mathcal{U}})}\right\} \cdot \mathscr{E}, \tag{16}$$

$$\max_{\mathcal{U}} \left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq \mathscr{E}. \tag{17}$$

**Corollary 1.** *Under the assumptions of Theorem 1, for all $\mathcal{U} \subseteq \{1, \ldots, d\}$*

$$\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq \min\left(1, \ \mathscr{E} + 2\sqrt{S_{\mathcal{U}}}, \ \mathscr{E} + 2\sqrt{1 - S_{\mathcal{U}}}\right) \cdot \mathscr{E}. \tag{18}$$

*In particular, if for some $\mathcal{U} \subseteq \{1, \ldots, d\}$ the Sobol' index or the total-effect index $S_{\mathcal{U}} \in \{0, 1\}$, then*

$$\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq \mathscr{E}^2. \tag{19}$$

The following corollary gives the bound for the sum of errors of Sobol' indices for all $2^d$ different subgroups of variables (not valid for total-effects).

**Corollary 2.** *For any functions $f, \widehat{f} \in L^2(\mathscr{X}, \mu)$ such that $\mathbb{V}_\mu[f] > 0$, $\mathbb{V}_\mu[\widehat{f}] > 0$, it holds for corresponding Sobol' indices for $\mathcal{U} \subseteq \{1, \ldots, d\}$*

$$\sum_{\mathcal{U}} \left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq 2 \cdot \mathscr{E}, \tag{20}$$

$$\sum_{\mathcal{U}} \left(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right)^2 \leq 2 \cdot \mathscr{E}^2. \tag{21}$$

Corollary 1 allows us to propose a new method for quality control of estimates of Sobol' indices based on metamodels (see Algorithm 1). The method uses the estimate

$$\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}\right| \leq \min\left(1, \ \mathscr{E}_2 + 2\sqrt{\widehat{S}_{\mathcal{U}}}, \ \mathscr{E}_2 + 2\sqrt{1 - \widehat{S}_{\mathcal{U}}}\right) \cdot \mathscr{E}_2, \tag{22}$$

where $\mathscr{E}_2 \triangleq \|f - \widehat{f}\|_\mu \cdot \min\left\{\mathbb{V}_\mu^{-1/2}[f], \mathbb{V}_\mu^{-1/2}[\widehat{f}]\right\}$, which follows from the symmetry of Theorem 1 with respect to $f$ and $\widehat{f}$. The values on the right-hand side of (22), which cannot be calculated analytically, are replaced by sample estimates. Approximation error is estimated using hold-out validation.

For simplicity, it is assumed that a metamodel of type (8) is used, and there is no additional random noise in the responses ($\sigma^2 = 0$). However, the method is easily generalized to both arbitrary metamodels and the noisy case. Asymptotic computational complexity[6] of Algorithm 1 is $\mathcal{O}(n_t)$, $n_t$ — test sample size.

One can show that the error upper bounds in Theorem 1 are achievable, using

---

[6]Number of regressors $N$, dimension $d$, and calculation time of the response $f(\mathbf{x}_i)$ are assumed to be constant.

**Algorithm 1.** Estimation of errors of Sobol' indices / total-effect indices.

---

**Parameters:** test set size $n_t$; constructed approximation $\widehat{f} = \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N} \widehat{c}_{\boldsymbol{\alpha}} \Psi_{\boldsymbol{\alpha}}$.

1. Obtain a sample $\mathcal{D}_t = (\mathbf{x}_i \in \mathscr{X})_{i=1}^{n_t}$ from the distribution $\mu$ and the responses $f(\mathbf{x}_i)$ and $\widehat{f}(\mathbf{x}_i)$ for each $\mathbf{x}_i \in \mathcal{D}_t$.

2. $\mathscr{M} \leftarrow \frac{1}{n_t} \sum_{i=1}^{n_t} \left[ f(\mathbf{x}_i) - \widehat{f}(\mathbf{x}_i) \right]^2$, where all $\mathbf{x}_i \in \mathcal{D}_t$.

3. $\mathscr{V}_1 \leftarrow \frac{1}{n_t} \sum_{i=1}^{n_t} \left[ f(\mathbf{x}_i) - \frac{1}{n_t} \sum_{j=1}^{n_t} f(\mathbf{x}_j) \right]^2$ where all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_t$

4. $\mathscr{V}_2 \leftarrow \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N^+} \widehat{c}_{\boldsymbol{\alpha}}^2$.

5. $\widehat{\mathscr{E}_2} \leftarrow \sqrt{\mathscr{M} / \max(\mathscr{V}_1, \mathscr{V}_2)}$.

6. For each $\mathcal{U} \subseteq \{1, \ldots, d\}$:

   6.1. Calculate $\widehat{S}_{\mathcal{U}}$ from the expansion coefficients $\widehat{c}_{\boldsymbol{\alpha}}$.

   6.2. $\mathscr{Q}_{\mathcal{U}} \leftarrow \min \left( 1, \ \widehat{\mathscr{E}_2} + 2\sqrt{\widehat{S}_{\mathcal{U}}}, \ \widehat{\mathscr{E}_2} + 2\sqrt{1 - \widehat{S}_{\mathcal{U}}} \right) \cdot \widehat{\mathscr{E}_2}$.

**Output:** upper bounds for $\left| S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}} \right|$ for all $\mathcal{U} \subseteq \{1, \ldots, d\}$, $\{\mathscr{Q}_{\mathcal{U}}\}$.

---

two functions of the form

$$
\begin{aligned}
f(\mathbf{x}) &= c_1 \cdot \Psi_{(1,0,0,\dots)}(\mathbf{x}) + c_2 \cdot \Psi_{(0,1,0,\dots)}(\mathbf{x}), \\
\widehat{f}(\mathbf{x}) &= \widehat{c}_1 \cdot \Psi_{(1,0,0,\dots)}(\mathbf{x}) + \widehat{c}_2 \cdot \Psi_{(0,1,0,\dots)}(\mathbf{x}).
\end{aligned}
\tag{23}
$$

**Theorem 2.** *For any subset $\mathcal{U} \subseteq \{1, \ldots, d\}$ and any values $S_{\mathcal{U}}, \widehat{S}_{\mathcal{U}} \in (0, 1)$, there are $f, \widehat{f} \in L^2(\mathscr{X}, \mu)$ with Sobol' (total-effect) indices for $\mathbf{x}_{\mathcal{U}}$ variables equal to $S_{\mathcal{U}}$ and $\widehat{S}_{\mathcal{U}}$ correspondingly, such that*

$$
\left| S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}} \right| = \left\{ \sqrt{S_{\mathcal{U}}(1 - \widehat{S}_{\mathcal{U}})} + \sqrt{\widehat{S}_{\mathcal{U}}(1 - S_{\mathcal{U}})} \right\} \cdot \mathscr{E}.
\tag{24}
$$

**Theorem 3.** *For any $\varepsilon \in [0, 1]$ there are $f, \widehat{f} \in L^2(\mathscr{X}, \mu)$ such that for their Sobol' indices (and for their total-effect indices) it holds*

$$
\max_{\mathcal{U} \subseteq \{1, \ldots, d\}} \left| S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}} \right| = \varepsilon \cdot \mathscr{E}.
\tag{25}
$$

In general, the error bound in Theorem 1 can be overestimated, but it is achievable in principle. Thus, the problem of accuracy of Sobol' indices estimates is reduced to the assessment of approximation quality.

**In the third chapter**, non-asymptotic risk upper bounds for Sobol' indices estimates for random design are obtained.

**The first part** deals with the risk of the estimates $\widehat{S}^P$ and $\widehat{S}^{LS}$ associated with the two methods of calculating the expansion coefficients. All risk bounds in this chapter are obtained under the condition of random design:

**Condition 1** (of random design). *Suppose the design of experiments $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$ consists of independent and identically distributed random variables with the distribution $\mu$.*

Assume that there is some fixed deterministic learning procedure that constructs the approximation based on the given training set (10):

$$\mathcal{L}\colon \big(\mathcal{D},\ Y = f(\mathcal{D}) + \boldsymbol{\eta}\big) \to \widehat{f} \tag{26}$$

such that a fixed realization of the training sample always leads to the same approximation $\widehat{f} \in L^2(\mathscr{X}, \mu)$.

The following theorem establishes a relationship between the risk of Sobol' indices estimates and the quadratic risk of the approximation $\mathbb{E}\|f - \widehat{f}\|_\mu^2$, where the expectation is taken with respect to distributions of design and noise. In what follows, we will assume that $\mathbb{V}_\mu[f] > 0$. Additionally define Sobol' indices $\widehat{S}_\mathcal{U} = 2^{-d}$ if $\mathbb{V}_\mu[\widehat{f}] = 0$.

**Theorem 4.** *Let $\widehat{f}$ be an arbitrary approximation of $f$ constructed according to the procedure $\mathcal{L}$ satisfying (26). Assume that under Condition 1 of random design there exists $\mathbb{E}\|f - \widehat{f}\|_\mu^2 < \infty$. Then for corresponding Sobol' indices and total-effect indices of $f$ and $\widehat{f}$ for $\mathcal{U} \subseteq \{1, \ldots, d\}$*

$$\max_{\mathcal{U}} \mathbb{E}\big(S_\mathcal{U} - \widehat{S}_\mathcal{U}\big)^2 \ \leq\ \mathcal{R}^2, \tag{27}$$

$$\mathbb{E}\big|S_\mathcal{U} - \widehat{S}_\mathcal{U}\big| \ \leq\ \mathcal{R}\left(\mathcal{R} + 2\sqrt{S_\mathcal{U}}\right), \tag{28}$$

$$where\ \ \mathcal{R}^2 \ \triangleq\ \frac{\mathbb{E}\|f - \widehat{f}\|_\mu^2}{\mathbb{V}_\mu[f]}. \tag{29}$$

**Corollary 3.** *Under the assumptions of Theorem 4, it holds for the corresponding Sobol' indices[7] of functions $f$ and $\widehat{f}$ for $\mathcal{U} \subseteq \{1, \ldots, d\}$*

$$\mathbb{E}\Big[\sum_{\mathcal{U}} \big(S_\mathcal{U} - \widehat{S}_\mathcal{U}\big)^2\Big] \leq 2 \cdot \mathcal{R}^2. \tag{30}$$

Projection method. Consider now not a general metamodel, but Polynomial Chaos Approximation.

---

[7]Not valid for total-effect indices.

**Condition 2** (of boundedness). *We additionally require $f$ to be bounded on $\mathscr{X}$:*

$$\left| f(\mathbf{x}) \right| \leq L \quad for \ \mathbf{x} \in \mathscr{X}. \tag{31}$$

**Theorem 5.** *Under Condition 1 of random design and Condition 2 of boundedness, for corresponding Sobol' (total-effect) indices of functions $f$ and $\widehat{f}^P$ it holds for $\mathcal{U} \subseteq \{1, \dots, d\}$*

$$\max_{\mathcal{U}} \mathbb{E} \left( S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^P \right)^2 \ \leq \ \mathcal{R}_p^2, \tag{32}$$

$$\mathbb{E} \left| S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^P \right| \ \leq \ \mathcal{R}_p \left( \mathcal{R}_p + 2\sqrt{S_{\mathcal{U}}} \right), \tag{33}$$

$$where \ \mathcal{R}_p^2 \ \triangleq \ \frac{1}{\mathbb{V}_\mu[f]} \cdot \|e_N\|_\mu^2 + \frac{L^2 + \sigma^2}{\mathbb{V}_\mu[f]} \cdot \frac{N}{n}.$$

**Corollary 4.** *Under the assumptions of Theorem 5, suppose additionally $\lim_{N \to \infty} \|e_N\|_\mu = 0$. Let $N = N(n)$,*

$$\frac{N}{n} \to 0 \ \ and \ \ N \to \infty \ \ as \ \ n \to \infty,$$

*Then*

$$\mathbb{E} \left( S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^P \right)^2 \underset{n \to \infty}{\longrightarrow} 0.$$

Ordinary least squares.  Introduce the numerical characteristic often used in the random design setting to control the "stability" of LS estimates:

**Definition 3.** *For the orthonormal set of functions $\{\Psi_{\boldsymbol{\alpha}}, \ \boldsymbol{\alpha} \in \mathscr{L}_N\}$ that satisfies (7), and for some fixed sequence of truncation sets (11), define*

$$K_N \triangleq \sup_{\mathbf{x} \in \mathscr{X}} \left[ \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N} \Psi_{\boldsymbol{\alpha}}^2(\mathbf{x}) \right]. \tag{34}$$

Denote *the spectral norm* of matrix $A \in \mathbb{R}^{m \times p}$ as

$$|||A||| = \max_{\mathbf{z} \in \mathbb{R}^p : \ \|\mathbf{z}\| \neq 0} \frac{\|A\mathbf{z}\|}{\|\mathbf{z}\|}. \tag{35}$$

We need the following result obtained by Cohen [6].

**Lemma 1** (Cohen, 2013). *Under Condition 1 of random design, for $\delta \in (0, 1)$*

$$P \left\{ |||\Phi^T \Phi / n - I_N||| > \delta \right\} \leq 2N \cdot \exp \left[ -\frac{c_\delta \cdot n}{K_N} \right], \tag{36}$$

*where $c_\delta \triangleq (1 + \delta) \ln(1 + \delta) - \delta > 0$.*

Lemma 1 leads to the condition on the size of the training sample $n$ and the number of regressors $N$ that excludes the possibility of ill-conditioned normalized information matrix $\Phi^T\Phi/n$ with high probability.

**Condition 3** (of stability). *Let for some fixed $r > 0$ the relation of $N$ and $n$ satisfies*

$$K_N \leq \varkappa_r \cdot \frac{n}{\ln n}, \quad where \quad \varkappa_r = \frac{3 \cdot \ln(3/2) - 1}{2 + 2r}. \tag{37}$$

Under Condition 3 of stability, we have based on Lemma 1

$$P\left\{|||\Phi^T\Phi/n - I_N||| > 1/2\right\} \leq 2n^{-r}. \tag{38}$$

**Theorem 6.** *Under Condition 1 of random design and Condition 3 of stability for corresponding Sobol' indices and total-effect indices of $f$ and $\widehat{f}^{LS}$ it holds for $\mathcal{U} \subseteq \{1, \ldots, d\}$*

$$\max_{\mathcal{U}} \mathbb{E}\left(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\right)^2 \leq \mathcal{R}_{LS}^2 + 2n^{-r}, \tag{39}$$

$$\mathbb{E}\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\right| \leq \mathcal{R}_{LS}\left(\mathcal{R}_{LS} + 2\sqrt{S_{\mathcal{U}}}\right) + 2n^{-r}, \tag{40}$$

$$where \quad \mathcal{R}_{LS}^2 \triangleq \frac{1.2}{\mathbb{V}_\mu[f]} \cdot \|e_N\|_\mu^2 + \frac{4\sigma^2}{\mathbb{V}_\mu[f]} \cdot \frac{N}{n}.$$

**Corollary 5.** *Under the assumptions of Theorem 6 for the case of noiseless observations, i.e. $\sigma^2 = 0$, it holds*

$$\mathbb{E}\left(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\right)^2 \leq \frac{1.2}{\mathbb{V}_\mu[f]}\|e_N\|_\mu^2 + 2n^{-r}, \tag{41}$$

$$\mathbb{E}\left|S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\right| \leq \frac{1.2}{\mathbb{V}_\mu[f]}\|e_N\|_\mu^2 + \frac{2.2 \cdot S_{\mathcal{U}}^{1/2}}{\mathbb{V}_\mu^{1/2}[f]}\|e_N\|_\mu + 2n^{-r}. \tag{42}$$

**Corollary 6.** *Under the assumptions of Theorem 6, except Condition 3, suppose additionally that $\lim_{N\to\infty} \|e_N\|_\mu = 0$. Let $N = N(n)$,*

$$\frac{K_N \cdot \ln N}{n} \to 0 \quad and \quad N \to \infty \quad as \quad n \to \infty,$$

*Then*

$$\mathbb{E}\left(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\right)^2 \xrightarrow[n\to\infty]{} 0.$$

**In the second part**, we perform an asymptotic analysis of the obtained risk upper bounds with increasing sample size for approximations based on specific families of multidimensional polynomials. We will assume that the analyzed function $f$ is *p-smooth*.

Table 1: Asymptotic upper bounds for the quadratic risk of Sobol' and total-effect indices estimates, depending on sample size $n$, input dimension $d$, and smoothness $p$.

| Polynomials | **Legendre** | **Chebyshev** | **Trigonometric** |
|---|---|---|---|
| Distribution | $U([-1,1]^d)$ | $Arc([-1,1]^d)$ | $U([0,1]^d)$ |
| $\sigma^2 = 0$: $\quad \mathbb{E}\big(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\big)^2 \lesssim$ | $\left(\frac{n}{\ln n}\right)^{-p/d}$ | $\left(\frac{n}{\ln n}\right)^{-2p/d}$ | $\left(\frac{n}{\ln n}\right)^{-2p/d}$ |
| $\mathbb{E}\big(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{P}\big)^2 \lesssim$ | $n^{-\frac{2p}{2p+d}}$ | | |
| $\sigma^2 > 0$: $\quad \mathbb{E}\big(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{LS}\big)^2 \lesssim$ | $n^{-\frac{2p}{2p+d}}, \quad p/d > 1/2$ <br> $\left(\frac{n}{\ln n}\right)^{-p/d}, \; p/d \le 1/2$ | $n^{-\frac{2p}{2p+d}}$ | $n^{-\frac{2p}{2p+d}}$ |
| $\mathbb{E}\big(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}}^{P}\big)^2 \lesssim$ | $n^{-\frac{2p}{2p+d}}$ | | |

As a basis for approximations, we consider three families of polynomials: Legendre, Trigonometric and Chebyshev, the first two of which are orthogonal with respect to the continuous uniform distribution, and the last — with respect to the arcsine distribution. In accordance with (7), the non-constant elements of these families are additionally normalized to have unit variance with respect to the corresponding distributions.

**Remark 1.** *For the case of Trigonometric polynomials, we additionally require that the function $f$ can be extended outside $\mathscr{X} = [0,1]^d$ to become 1-periodic in each input argument.*

For asymptotic analysis, *the truncation scheme* will be used based on the maximum degree of one-dimensional polynomial factors. For some $\alpha_{max} \in \mathbb{N}_+$ we define

$$\mathscr{L}_N = \{\boldsymbol{\alpha} \in \mathbb{N}^d \colon \max_{i=1,\dots,d}\{\alpha_i\} \le \alpha_{max}\}, \tag{43}$$

where $N = |\mathscr{L}_N| = (\alpha_{max} + 1)^d$.

Table 1, based on the results of Theorems 5 and 6, summarizes the asymptotic[8] risk bounds for estimates of Sobol' indices and total-effect indices for the two methods for calculating the expansion coefficients and three types of polynomials. When deriving these bounds, it was assumed that the number of regressors $N$ is chosen asymptotically optimally to minimize the resulting risk.

---

[8]If two sequences of positive numbers $\{a_n\}$ and $\{b_n\}$ are given, then $a_n \lesssim b_n$ means that the sequence $\{a_n/b_n\}$ is bounded.

It can be concluded that the key factors that provide the possibility for fast convergence of Sobol' indices estimates are the absence of random noise in the output of the analyzed function, its high smoothness, and low dimension.

**In the fourth chapter**, we consider another – asymptotic – approach for analyzing the quality of estimates of Sobol' indices and propose a method of sequential design of experiments based on it, which ensures high accuracy of these estimates.

In this chapter, we only consider Sobol' indices of the type $S_i \triangleq S_{\{i\}}$, $i = 1, \ldots, d$, called *first-order sensitivity indices*. In addition, we use a simplified[9] data model in which DoE is fixed (not random), and the analyzed function has the form of Polynomial Chaos Approximation (8) with a finite number of terms $f(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N} c_{\boldsymbol{\alpha}} \Psi_{\boldsymbol{\alpha}}(\mathbf{x})$, where $\mathbb{V}_\mu[f] = \sum_{\boldsymbol{\alpha} \in \mathscr{L}_N^+} c_{\boldsymbol{\alpha}}^2 > 0$. The response of the model is formed as

$$y = \mathbf{c}^T \boldsymbol{\Psi}(\mathbf{x}) + \eta, \tag{44}$$

where $\eta \sim N(0, \sigma^2)$, $\sigma^2 > 0$ is independent *Gaussian* noise, and $\boldsymbol{\Psi}(\mathbf{x}) \triangleq \big(\Psi_{\boldsymbol{\alpha}}(\mathbf{x}), \boldsymbol{\alpha} \in \mathscr{L}_N\big)^T \in \mathbb{R}^N$, and the metamodel has the form $\widehat{f}(\mathbf{x}) = \widehat{\mathbf{c}}^T \boldsymbol{\Psi}(\mathbf{x})$. The expansion coefficients are estimated from the training sample with the help of LS method only.

Define *information matrix* $A_n \in \mathbb{R}^{N \times N}$ as

$$A_n = \sum_{i=1}^{n} \boldsymbol{\Psi}(\mathbf{x}_i) \boldsymbol{\Psi}^T(\mathbf{x}_i). \tag{45}$$

Define a vector-valued function of first-order sensitivity indices with components of the form (9) as $\mathbf{S}(\mathbf{c}) \triangleq \big(S_1(\mathbf{c}), \ldots, S_d(\mathbf{c})\big)^T$ with the corresponding Jacobian matrix (of size $d \times N$)

$$B \triangleq B(\mathbf{c}) = \frac{\partial \mathbf{S}(\mathbf{c})}{\partial \mathbf{c}}. \tag{46}$$

**Theorem 7.**
*Let the following conditions be satisfied:*

1. *There is some infinite deterministic sequence of points in the design space $\{\mathbf{x}_i \in \mathscr{X}\}_{i=1}^{\infty}$ such that*

$$\frac{1}{n} A_n \xrightarrow[n \to \infty]{} H, \tag{47}$$

   *where $H \in \mathbb{R}^{N \times N}$ — some symmetric positive definite matrix.*

2. *The points from this sequence $\{\mathbf{x}_i \in \mathscr{X}\}_{i=1}^{\infty}$ and the corresponding responses (44) are iteratively added to the training set.*

3. *For the true coefficients $\mathbf{c}$, $\det\big(B(\mathbf{c}) H^{-1} B^T(\mathbf{c})\big) \neq 0$ holds.*

---

[9]Thus, the difference from the previously presented data model is that here $e_N(\mathbf{x}) \equiv 0$, the noise is Gaussian, and Condition 1 of random design is not imposed. In addition, the structure of the metamodel is fixed.

---

**Algorithm 2.** DoE for estimating Sobol' indices.

---

**Parameters:** initial number of points in the design of experiments $m$ and the final number $n > m$; set of candidate points $\Xi$.

**Initialization:** initial training set $(\mathcal{D}, Y)$ of $m$ examples such that the design $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^m \subset \Xi$ leads to a nonsingular information matrix $A_m = \boldsymbol{\Psi}^T(\mathcal{D}) \cdot \boldsymbol{\Psi}(\mathcal{D})$; estimates $\widehat{\mathbf{c}}_m = A_m^{-1} T_m$, where $T_m = \boldsymbol{\Psi}^T(\mathcal{D}) \, Y$.

**Iterations:** while the design $\mathcal{D}$ contains less than $n$ examples:

1. $\mathbf{x} \leftarrow \arg\min_{\mathbf{x} \in \Xi} \det\left[ B(\widehat{\mathbf{c}}) \cdot \left\{ A + \boldsymbol{\Psi}(\mathbf{x}) \cdot \boldsymbol{\Psi}^T(\mathbf{x}) \right\}^{-1} \cdot B^T(\widehat{\mathbf{c}}) \right].$     // see (49)

2. $\mathcal{D} \leftarrow \mathcal{D} \cup \{\mathbf{x}\}, \quad Y \leftarrow Y \cup \{f(\mathbf{x})\}.$

3. $A \leftarrow A + \boldsymbol{\Psi}(\mathbf{x}) \cdot \boldsymbol{\Psi}^T(\mathbf{x}), \quad T \leftarrow T + f(\mathbf{x}) \cdot \boldsymbol{\Psi}(\mathbf{x}).$

4. $\widehat{\mathbf{c}} \leftarrow A^{-1} T.$         // update estimates according to (14)

**Output:** design of experiments, $\mathcal{D}$.

---

*Then the following convergence in distribution takes place:*

$$\sqrt{n}\big(\mathbf{S}(\mathbf{c}) - \mathbf{S}(\widehat{\mathbf{c}}_n)\big) \xrightarrow[n\to\infty]{d} N\big(0, \, \sigma^2 B H^{-1} B^T\big). \tag{48}$$

Theorem 7 leads to the idea of Algorithm 2 of sequential DoE for estimating Sobol' indices. The algorithm uses the $D$-optimality criterion, and at each iteration of the algorithm, the current estimates of the expansion coefficients and the current normalized information matrix are used to calculate the determinant of the covariance matrix $\sigma^2 \cdot B(\mathbf{c}) H^{-1} B^T(\mathbf{c})$. Instead of calculating the determinant for each candidate point at step 1. of Algorithm 2, we use an equivalent optimization problem with a computationally efficient solution:

$$\frac{\boldsymbol{\Psi}^T(\mathbf{x}) \cdot \left\{ A^{-1} B^T (B A^{-1} B^T)^{-1} B A^{-1} \right\} \cdot \boldsymbol{\Psi}(\mathbf{x})}{1 + \boldsymbol{\Psi}^T(\mathbf{x}) \cdot A^{-1} \cdot \boldsymbol{\Psi}(\mathbf{x})} \to \max_{\mathbf{x} \in \Xi}. \tag{49}$$

Asymptotic computational complexity[10] of Algorithm 2 is $\mathcal{O}(n)$, where $n$ is the final size of the design after adding all new points.

In general, our approach to DoE is to obtain an (asymptotic) normal distribution of Sobol' indices with a design-dependent covariance matrix and apply one of the optimality criteria. Algorithm 2 illustrates this idea using the $D$-optimality criterion, but

---

[10]Number of regressors $N$, dimension $d$, number of candidate points $|\Xi|$, and calculation time of the response $f(\mathbf{x}_i)$ are assumed to be constant.

other criteria can be used. For example, Pronzato [12] applied the proposed approach using the $A$- and $MV$-criteria. The choice of a specific criterion for the design optimality should be carried out by the researcher based on the specifics of the problem. In particular, the $D$-optimality criterion can be recommended when the "mean" error of the estimates of Sobol' indices is more important than its maximum value over all groups of input parameters.

**The fifth chapter** describes the software package developed by the author and gives the results of computational experiments for the proposed algorithms.

**The first part** is devoted to the developed software package in the Python language, which includes the algorithms for quality control and design of experiments created in the study. In addition to the algorithms, the package also includes a test environment (test analyzed functions and alternative methods from the literature), which allows you to compare the proposed approaches with analogues.

**In the second part**, we test the quality control method for Sobol' indices (Algorithm 1). For comparison, we use sample error bounds based on the *bootstrap* method [7]. An example of results for the 2-dimensional Sobol' $g$-function is shown in Figure 1a. When assessing the quality of the approximation, 15% of the samples are used for hold-out validation.

**In the third part**, an experimental analysis of the risk bounds obtained in Theorems 5 and 6 is given, and it is shown that the metamodeling approach does allow one to achieve a high rate of convergence of estimates of Sobol' indices to their true values.
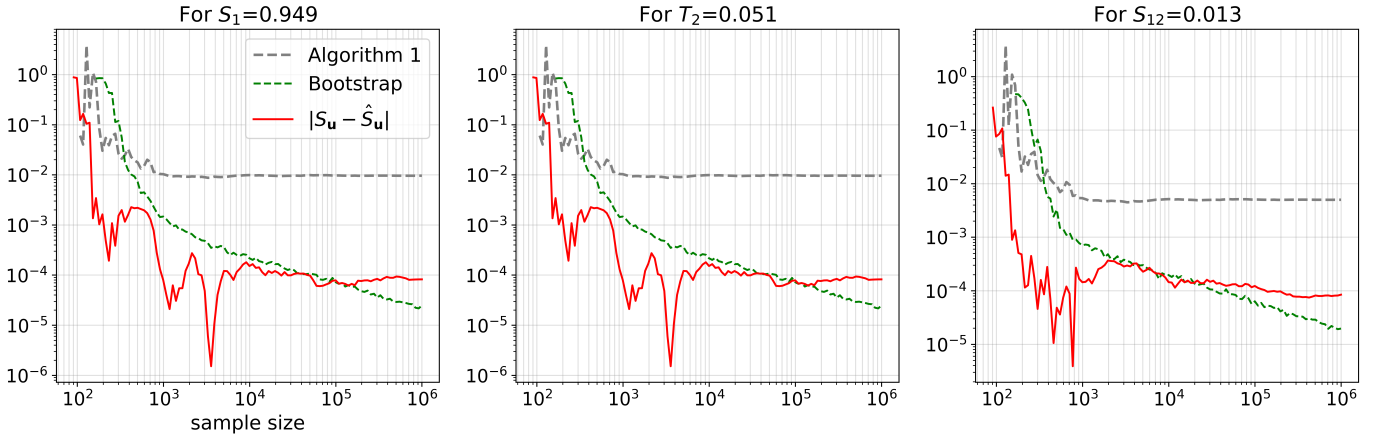
As an illustration, Figure 1b shows the empirical estimate of risk $\max_{\mathcal{U}} \left\{ \mathbb{E}(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}})^2, \ \mathbb{E}(T_{\mathcal{U}} - \widehat{T}_{\mathcal{U}})^2 \right\}$ and the components of its bounds from the mentioned theorems $\|e_N\|_{\mu}^2 / \mathbb{V}_{\mu}[f]$ and $n^{-r}$ for Sobol' $g$-function in the noiseless case.

**The fourth part** presents the results of applying the proposed method of sequential design of experiments (Algorithm 2) to solve a series of artificial and real engineering problems (using finite element models), the dimension of the design space varies from 2 to 53. The experiment setting assumes that new points are iteratively added to some initial random design. The proposed method is compared with the following DoE techniques:
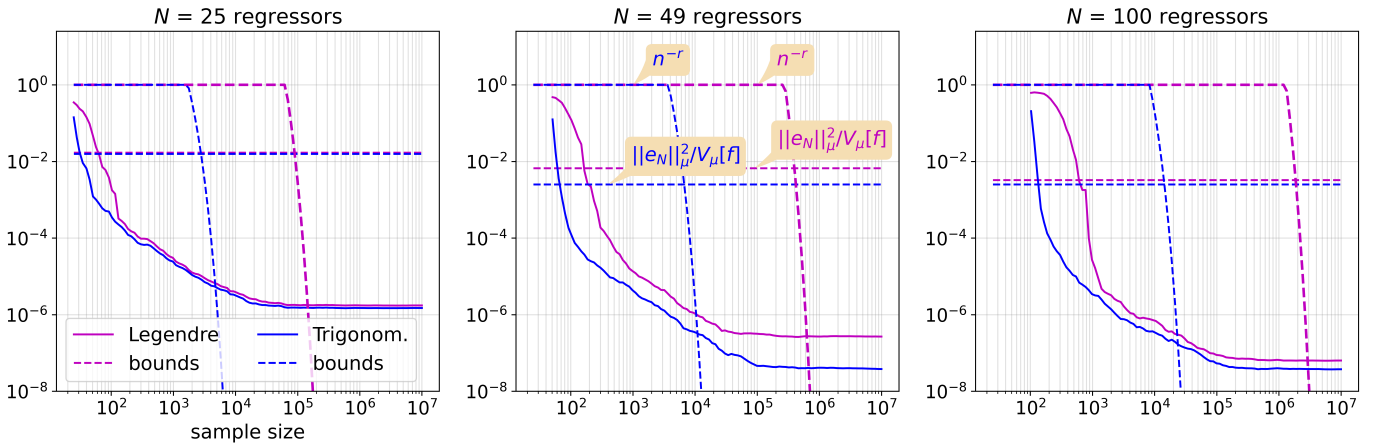
- *Random design*, the sequential addition of random points from a set of candidates $\Xi$.

- *Sequential D-optimal design* [13] based on iterative maximization of determinant of the information matrix: $\mathbf{x}_{n+1} = \arg\max_{\mathbf{x} \in \Xi} \det \left[ A_n + \boldsymbol{\Psi}(\mathbf{x}) \cdot \boldsymbol{\Psi}^T(\mathbf{x}) \right]$.

- *LHS* — sampling based on *latin hypercube*. Note that at each iteration all points of the design are updated.

The effectiveness of the proposed approach is illustrated by Figure 1c, which shows the results of the listed DoE techniques for the model of deflection of a bar structure (truss) under the action of external forces [2, 9], in which $d = 10$ input parameters have a continuous uniform distribution. The metric of DoE quality is the error of
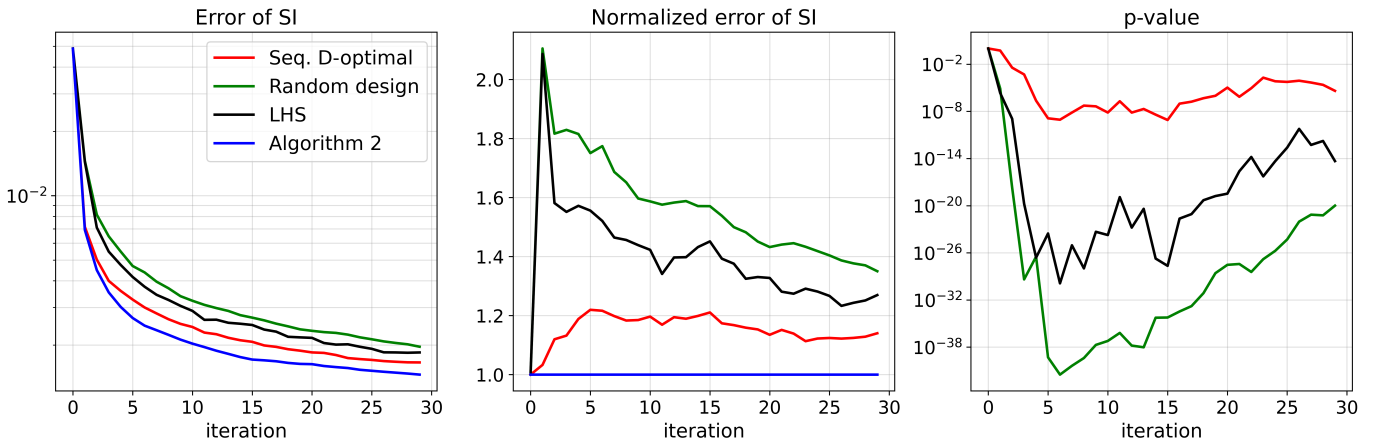
Sobol' indices $\sqrt{\sum_{i=1}^{d}(S_i - \widehat{S}_i)^2}$ averaged over several runs of the DoE technique with different initial designs. For convenience, the figure also shows the error values for all techniques relative to Algorithm 2. In addition, we use Welch's $t$-test to verify that the difference in the mean errors of the indices is statistically significant. It can be seen that the estimates of Sobol' indices obtained on the basis of Algorithm 2 are, on average, more accurate at all iterations; in particular, after adding 29 new points, the average accuracy of these estimates by 10% exceeds the result that gives the most accurate of the other methods in comparison, the sequential $D$-optimal design.

(a) Estimation errors of the three Sobol' indices and their sample-based error bounds for Algorithm 1 and the bootstrap method. Metamodels based on Legendre polynomials, $N = 91$ regressors, LS.



(b) Quadratic risk of Sobol' and total-effect indices $\max_{\mathcal{U}} \left\{ \mathbb{E}(S_{\mathcal{U}} - \widehat{S}_{\mathcal{U}})^2, \ \mathbb{E}(T_{\mathcal{U}} - \widehat{T}_{\mathcal{U}})^2 \right\}$ and the components of its bounds according to Theorem 6. Metamodels based on Legendre and Trigonometric polynomials with different number of regressors. LS. Noiseless case, $\sigma^2 = 0$.



(c) The mean error of Sobol' indices estimates and its normalized version; and $p$-value of Welch's $t$-test for different DoE. Legendre polynomials, $N = 176$ regressors, LS. Initial design size $m = 176$.

Figure 1: Some results of computational experiments.

# Conclusion

1. A relationship has been established between the error in estimates of Sobol' indices and the error of the approximation, on the basis of which these estimates were obtained. This relationship is valid for Sobol' indices and total-effect indices of all orders. In particular, it is shown that the maximum absolute error of estimates of Sobol' indices for all groups of variables is bounded by the relative error of the corresponding approximation, and this bound is achievable.

2. Thanks to the obtained theoretical error bound, a method for controlling the quality of metamodel-based estimates of Sobol' indices has been developed.

3. Under the condition of random design of experiments for Polynomial Chaos Approximation, non-asymptotic upper bounds on the risk of metamodel-based estimates of Sobol' indices are obtained. In addition, estimates for the rate of convergence are found for these bounds in the case of analyzed functions of different smoothness and approximations that use Legendre, Chebyshev and Trigonometric polynomials.

4. An asymptotic distribution of estimates of Sobol' indices was found, which made it possible to develop a method for sequential design of experiments for estimating sensitivity indices using Polynomial Chaos Approximation.

5. A software package has been developed for solving problems related to modeling in engineering design, which includes the proposed methods for quality control of estimates of Sobol' indices and the design of experiments.

6. The effectiveness of the developed software package was demonstrated in solving a number of engineering problems; particularly, in the analysis of factors that affect the deflection of a bar structure (truss) under the action of external forces.

# References

[1] Jian An and Art Owen. "Quasi-regression". In: vol. 17. 4. Complexity of multivariate problems (Kowloon, 1999). 2001, pp. 588–607. DOI: `10.1006/jcom.2001.0588`. URL: `https://doi.org/10.1006/jcom.2001.0588`.

[2] Géraud Blatman and Bruno Sudret. "An adaptive algorithm to build up sparse polynomial chaos expansions for stochastic finite element analysis". In: *Probabilistic Engineering Mechanics* 25 (Apr. 2010), pp. 183–197. DOI: `10.1016/j.probengmech.2009.10.003`.

[3] Evgeny Burnaev and Ivan Panin. "Adaptive Design of Experiments for Sobol Indices Estimation Based on Quadratic Metamodel". In: *Statistical Learning and Data Sciences*. Ed. by Alexander Gammerman, Vladimir Vovk, and Harris Papadopoulos. Cham: Springer International Publishing, 2015, pp. 86–95. ISBN: 978-3-319-17091-6.

[4] Evgeny Burnaev, Ivan Panin, and Bruno Sudret. "Effective Design for Sobol Indices Estimation Based on Polynomial Chaos Expansions". In: *Conformal and Probabilistic Prediction with Applications*. Ed. by Alexander Gammerman et al. Cham: Springer International Publishing, 2016, pp. 165–184. ISBN: 978-3-319-33395-3.

[5] Evgeny Burnaev, Ivan Panin, and Bruno Sudret. "Efficient design of experiments for sensitivity analysis based on polynomial chaos expansions". In: *Ann. Math. Artif. Intell.* 81.1-2 (2017), pp. 187–207. ISSN: 1012-2443. DOI: `10.1007/s10472-017-9542-1`. URL: `https://doi.org/10.1007/s10472-017-9542-1`.

[6] Albert Cohen, Mark A. Davenport, and Dany Leviatan. "On the stability and accuracy of least squares approximations". In: *Found. Comput. Math.* 13.5 (2013), pp. 819–834. ISSN: 1615-3375. DOI: `10.1007/s10208-013-9142-3`. URL: `https://doi.org/10.1007/s10208-013-9142-3`.

[7] S. Dubreuil et al. "Construction of bootstrap confidence intervals on sensitivity indices computed by polynomial chaos expansion". In: *Reliability Engineering & System Safety* 121 (2014), pp. 263–275. ISSN: 0951-8320. DOI: `https://doi.org/10.1016/j.ress.2013.09.011`. URL: `http://www.sciencedirect.com/science/article/pii/S0951832013002688`.

[8] Bertrand Iooss and Paul Lemaître. "A Review on Global Sensitivity Analysis Methods". In: *Operations Research/ Computer Science Interfaces Series* 59 (Apr. 2014). DOI: `10.1007/978-1-4899-7547-8_5`.

[9] Sang Hoon Lee and Byung Man Kwak. "Response surface augmented moment method for efficient reliability analysis". In: *Structural Safety* 28.3 (2006), pp. 261–272. ISSN: 0167-4730. DOI: `https://doi.org/10.1016/j.strusafe.2005.08.003`. URL: `https://www.sciencedirect.com/science/article/pii/S0167473005000421`.

[10] Ivan Panin. "Risk of estimators for Sobol' sensitivity indices based on metamodels". In: *Electron. J. Statist.* 15.1 (2021), pp. 235–281. ISSN: 1935-7524. DOI: 10.1214/20-EJS1793. URL: https://projecteuclid.org/euclid.ejs/1609902190.

[11] Ivan Panin and Pavel Prikhodko. "Approaches to the evaluation of the sensitivity indices variance in the problem of global sensitivity analysis". Russian. In: *Proceedings of the conference "Information Technology and Systems"* (Petrozavodsk, Russia). IITP RAS. 2012, pp. 173–178. ISBN: 978-5-901158-19-7. URL: http://www.itas2012.iitp.ru/pdf/1569602539.pdf.

[12] Luc Pronzato. "Sensitivity analysis via Karhunen-Loève expansion of a random field model: Estimation of Sobol' indices and experimental design". In: *Reliability Engineering & System Safety* (Jan. 2018). DOI: 10.1016/j.ress.2018.01.010.

[13] Luc Pronzato and Andrej Pázman. *Design of experiments in nonlinear models.* Vol. 212. Lecture Notes in Statistics. Asymptotic normality, optimality criteria and small-sample properties. Springer, New York, 2013, pp. xvi+399. ISBN: 978-1-4614-6362-7; 978-1-4614-6363-4. DOI: 10.1007/978-1-4614-6363-4. URL: https://doi.org/10.1007/978-1-4614-6363-4.

[14] Andrea Saltelli et al. *Global sensitivity analysis. The primer.* John Wiley & Sons, Ltd., Chichester, 2008, pp. xii+292. ISBN: 978-0-470-05997-5.

[15] I. M. Sobol'. "Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates". In: vol. 55. 1-3. The Second IMACS Seminar on Monte Carlo Methods (Varna, 1999). 2001, pp. 271–280. DOI: 10.1016/S0378-4754(00)00270-6. URL: https://doi.org/10.1016/S0378-4754(00)00270-6.

[16] I. M. Sobol'. "Sensitivity estimates for nonlinear mathematical models". In: *Math. Modeling Comput. Experiment* 1.4 (1993), 407–414 (1995). ISSN: 1061-7590.

[17] Bruno Sudret. "Polynomial chaos expansions and stochastic finite element methods". In: *Risk and Reliability in Geotechnical Engineering.* Ed. by Jianye Ching Kok-Kwang Phoon. CRC Press, 2015, pp. 265–300. URL: https://hal.archives-ouvertes.fr/hal-01449883.

[18] Bruno Sudret. "Uncertainty propagation and sensitivity analysis in mechanical models – Contributions to structural reliability and stochastic spectral methods". In: (Jan. 2007).

[19] A.W. van der Vaart. *Asymptotic Statistics.* Asymptotic Statistics. Cambridge University Press, 2000. ISBN: 9780521784504. URL: https://books.google.ru/books?id=UEuQEM5RjWgC.