

Sensitivity to statistical estimation uncertainties and probabilistic model identification

ETICS 2022

Charles Surget^{1,2}

charles.surget@onera.fr

Sylvain Dubreuil¹, Jérôme Morio¹, Cécile Mattrand², Jean-Marc Bourinet², Nicolas Gayton².

7th of October 2022

¹ ONERA/DTIS, F-31055 Toulouse, France ² SIGMA Clermont, F-63000 Clermont-Ferrand, France

Ce document est la propriété de l'ONERA. Il ne peut être communiqué à des tiers et/ou reproduit sans l'autorisation préalable écrite de l'ONERA, et son contenu ne peut être divulgué.

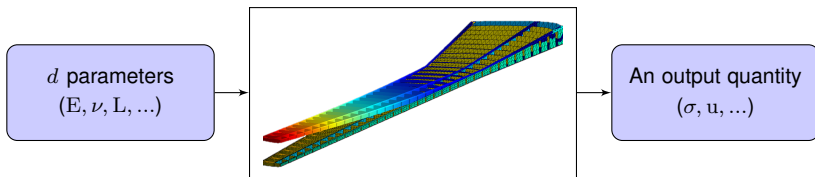
This document and the information contained herein is proprietary information of ONERA and shall not be disclosed or reproduced without the prior authorization of ONERA.

Table of contents

- 1 Context
- 2 Proposed approach
- 3 Sensitivity analysis
- 4 Illustration
- 5 Conclusion

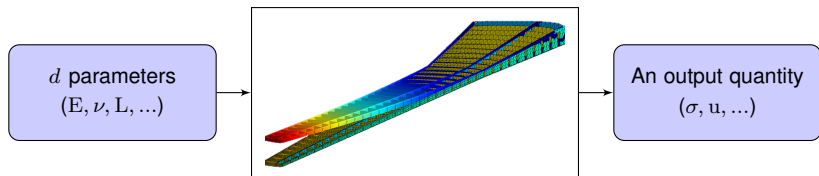
Context - Uncertainty Quantification

Simulation code of a mechanical structure:

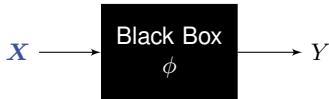


Context - Uncertainty Quantification

Simulation code of a mechanical structure:



In an uncertainty quantification context, those parameters are considered as an input continuous random vector :



with $X = (X_1, \dots, X_d)^t$ with values on the domain $\mathcal{X} \subseteq \mathbb{R}^d$ and defined by a given Probability Density Function (PDF) f_X .

Context - 1st uncertainty source

One could be interested in assessing the following expectation of a particular function τ of $Y = \phi(\mathbf{X})$ (e.g. a mean or a probability of failure):

$$\mathbb{E}_{f_{\mathbf{X}}} [\tau(\phi(\mathbf{X}))] = \int_{\mathcal{X}} \tau(\phi(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}. \quad (1.1)$$

Context - 1st uncertainty source

One could be interested in assessing the following expectation of a particular function τ of $Y = \phi(\mathbf{X})$ (e.g. a mean or a probability of failure):

$$\mathbb{E}_{f_{\mathbf{X}}} [\tau(\phi(\mathbf{X}))] = \int_{\mathcal{X}} \tau(\phi(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}. \quad (1.1)$$

Assuming $\tau = \text{Id}$, the Monte Carlo (MC) estimator of this integral is given by:

$$\hat{\mu}^{MC} = \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^{N_{\mathbf{X}}} \phi(\mathbf{X}^{(j)}), \quad (1.2)$$

with $\mathbf{X}^{(j)} \stackrel{i.i.d.}{\sim} f_{\mathbf{X}}$ and $N_{\mathbf{X}}$ the size of the MC sample.

Context - 1st uncertainty source

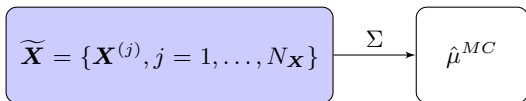
One could be interested in assessing the following expectation of a particular function τ of $Y = \phi(\mathbf{X})$ (e.g. a mean or a probability of failure):

$$\mathbb{E}_{f_{\mathbf{X}}} [\tau(\phi(\mathbf{X}))] = \int_{\mathcal{X}} \tau(\phi(\mathbf{x})) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}. \quad (1.1)$$

Assuming $\tau = \text{Id}$, the Monte Carlo (MC) estimator of this integral is given by:

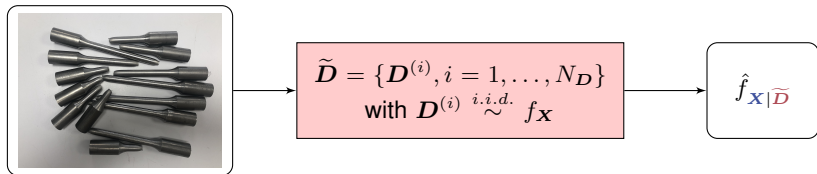
$$\hat{\mu}^{MC} = \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^{N_{\mathbf{X}}} \phi(\mathbf{X}^{(j)}), \quad (1.2)$$

with $\mathbf{X}^{(j)} \stackrel{i.i.d.}{\sim} f_{\mathbf{X}}$ and $N_{\mathbf{X}}$ the size of the MC sample. A first uncertainty source is related to this sample, defined as $\widetilde{\mathbf{X}}$ in the following process:



Context - 2nd uncertainty source

In a realistic context, the PDF f_X may be unknown [1]. Thus, the probabilistic model must be inferred from experimental tests:



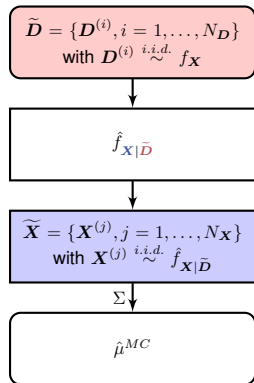
with N_D the size of the **database** \tilde{D} . The estimation $\hat{f}_{X|\tilde{D}}$ [2, 3] of the PDF f_X induces a second uncertainty source related to \tilde{D} .

[1] G Sarazin. Analyse de sensibilité fiabiliste en présence d'incertitudes épistémiques introduites par les données d'apprentissage, 2021.

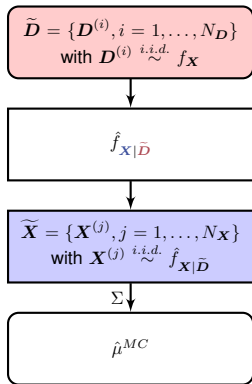
[2] K James, Lindsey and others. Parametric statistical inference. Oxford University Press, 1996.

[3] A J Izenman. Recent developments in nonparametric density estimation. Journal of the american statistical association, 1991.

Context - Small-Data

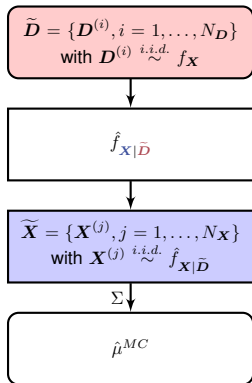


Context - Small-Data



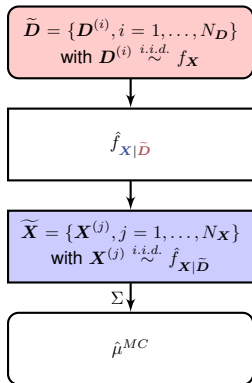
- The **database** is of limited size N_D : the small-data context is imposed by costly experimental tests.

Context - Small-Data



- The **database** is of limited size N_D : the small-data context is imposed by costly experimental tests.
- The **MC sample** is of limited size N_X : the small-data context is imposed by the simulation time induced by the model complexity.

Context - Small-Data



- The **database** is of limited size N_D : the small-data context is imposed by costly experimental tests.
- The **MC sample** is of limited size N_X : the small-data context is imposed by the simulation time induced by the model complexity.

↓

Test-Simulation trade-off

Context - Conditioning on a given database

The expectation (1.1) and its estimator are thus written as following for a given database $\tilde{D} = \tilde{d}$:

$$\mathbb{E}_{\hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}}[\phi(\mathbf{X})] = \int_{\mathcal{X}} \phi(\mathbf{x}) \hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}(\mathbf{x}) d\mathbf{x} \quad (1.3)$$

$$\approx \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^{N_{\mathbf{X}}} \phi(\mathbf{X}^{(j)}), \quad (1.4)$$

with $\mathbf{X}^{(j)} \stackrel{i.i.d.}{\sim} \hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}$. The estimator (1.4) is now subject to the first uncertainty source conditioned on the database $\tilde{D} = \tilde{d}$.

Context - Conditioning on a given database

The expectation (1.1) and its estimator are thus written as following for a given database $\tilde{D} = \tilde{d}$:

$$\mathbb{E}_{\hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}}[\phi(\mathbf{X})] = \int_{\mathbf{X}} \phi(\mathbf{x}) \hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}(\mathbf{x}) d\mathbf{x} \quad (1.3)$$

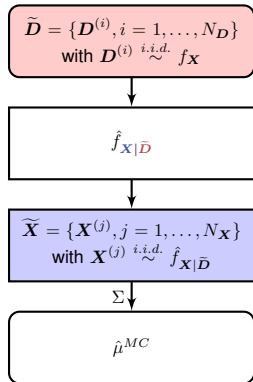
$$\approx \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^{N_{\mathbf{X}}} \phi(\mathbf{X}^{(j)}), \quad (1.4)$$

with $\mathbf{X}^{(j)} \stackrel{i.i.d.}{\sim} \hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}$. The estimator (1.4) is now subject to the first uncertainty source conditioned on the database $\tilde{D} = \tilde{d}$.

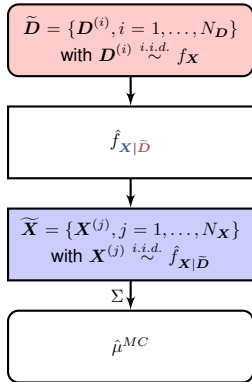
However, the uncertainty related to the database is not considered in the following variance :

$$\mathbb{V}_{\hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}}[\hat{\mu}^{MC}] = \frac{1}{N_{\mathbf{X}}} \mathbb{V}_{\hat{f}_{\mathbf{X}|\tilde{D}=\tilde{d}}}[\phi(\mathbf{X})]. \quad (1.5)$$

Context - Problems



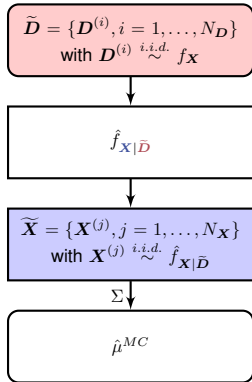
Context - Problems



Problem A

How to take into account the uncertainty of the **database** in the variance of the estimator?

Context - Problems



Problem A

How to take into account the uncertainty of the **database** in the variance of the estimator?

Problem B

In order to improve efficiently the accuracy of the estimator, should the investment of data be made in the **database** or the **MC** sample?

Table of contents

- 1 Context
- 2 Proposed approach
- 3 Sensitivity analysis
- 4 Illustration
- 5 Conclusion

Proposed approach - Double integral expectation

Problem A

How to take into account the uncertainty of the **database** in the variance of the estimator?

Proposed approach - Double integral expectation

Problem A

How to take into account the uncertainty of the **database** in the variance of the estimator?

The following expectation takes into account the variation of the **database**:

$$\mathbb{E}_{f_{(\mathbf{X}, \tilde{\mathbf{D}})}} [\phi(\mathbf{X})] = \int_{\mathbf{x}^{N_D}} \int_{\mathbf{x}} \phi(\mathbf{x}) f_{(\mathbf{X}, \tilde{\mathbf{D}})}(\mathbf{x}, \tilde{\mathbf{d}}) d\mathbf{x} d\tilde{\mathbf{d}}. \quad (2.1)$$

Proposed approach - Double integral expectation

Problem A

How to take into account the uncertainty of the **database** in the variance of the estimator?

The following expectation takes into account the variation of the **database**:

$$\mathbb{E}_{f(\mathbf{x}, \tilde{\mathbf{D}})} [\phi(\mathbf{X})] = \int_{\mathbf{x}^{N_D}} \int_{\mathbf{x}} \phi(\mathbf{x}) f_{(\mathbf{X}, \tilde{\mathbf{D}})}(\mathbf{x}, \tilde{\mathbf{d}}) d\mathbf{x} d\tilde{\mathbf{d}}. \quad (2.1)$$

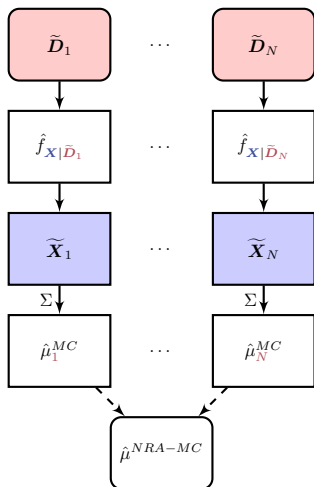
An estimator of this integral [4] is the following:

$$\hat{\mu}^{NRA-MC} = \frac{1}{N} \sum_{k=1}^N \frac{1}{N_X} \sum_{j=1}^{N_X} \phi(\mathbf{X}_k^{(j)}) = \frac{1}{N} \sum_{k=1}^N \hat{\mu}_k^{MC}, \quad (2.2)$$

with $\mathbf{X}_k^{(j)}$ *i.i.d.* $\hat{f}_{\mathbf{X}|\tilde{\mathbf{D}}_k}$ and N the number of **databases** of size N_D .

[4] V Chabridon. Analyse de sensibilité fiabiliste avec prise en compte d'incertitudes sur le modèle probabiliste, 2018.

Proposed approach - Empirical variance



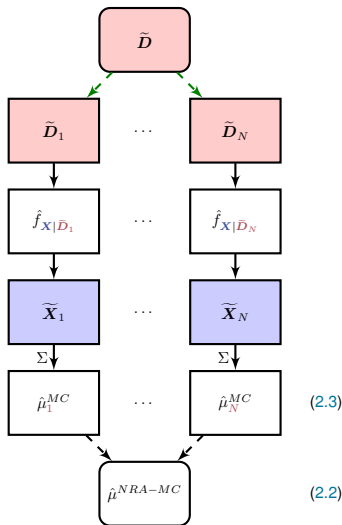
The process is repeated in order to estimate the variance with respect to the **database**.

Empirical variance:

$$\hat{V}_{f_{(x,\bar{D})}}[\hat{\mu}^{MC}] = \frac{1}{N-1} \sum_{k=1}^N (\hat{\mu}_k^{MC} - \hat{\mu}^{NRA-MC})^2 \quad (2.3)$$

(2.2)

Proposed approach - Small-data context



In a small-data context, only a **database \tilde{D}** of limited size N_D is available.

Resampling method

Allows to generate N **databases** from the initial one. [5, 6] (e.g. Bootstrap)

Solution A

The nested estimator (2.2) is now conditioned on the initial **database** but the method allows to take into account the uncertainty related to it.

[5] C H Yu. Resampling methods: concepts, applications, and justification. Practical Assessment, Research, and Evaluation, 2002.

[6] B Efron. The jackknife, the bootstrap and other resampling plans. SIAM, 1982.

Table of contents

- 1 Context
- 2 Proposed approach
- 3 Sensitivity analysis**
- 4 Illustration
- 5 Conclusion

Sensitivity analysis - ANOVA

Problem B

In order to improve efficiently the accuracy of the estimator, should the investment of data be made in the **database** or the **MC sample**?

Sensitivity analysis - ANOVA

Problem B

In order to improve efficiently the accuracy of the estimator, should the investment of data be made in the **database** or the **MC sample**?

An ANalysis Of VAriance [7, 8] is performed:

$$\left\{ \begin{array}{l} S_{\tilde{D}} = \frac{\mathbb{V} \left[\mathbb{E} \left[\hat{\mu}^{MC} \mid \tilde{D} \right] \right]}{\mathbb{V} \left[\hat{\mu}^{MC} \right]} \\ S_{\tilde{X}} = \frac{\mathbb{V} \left[\mathbb{E} \left[\hat{\mu}^{MC} \mid \tilde{X} \right] \right]}{\mathbb{V} \left[\hat{\mu}^{MC} \right]} \end{array} \right. \quad (3.1)$$

Interpretation of Sobol' indices:

- $S_{\tilde{D}}$ → proportion of variance due to the **database**,
- $S_{\tilde{X}}$ → proportion of variance due to the **MC sample**.

[7] Ilya M Sobol'. Sensitivity analysis for non-linear mathematical models. Mathematical modelling and computational experiment, 1993.

[8] F Gamboa and others. Statistical inference for sobol pick-freeze monte carlo method. Statistics, 50(4):881–902, 2016.

Sensitivity analysis - Interpretation

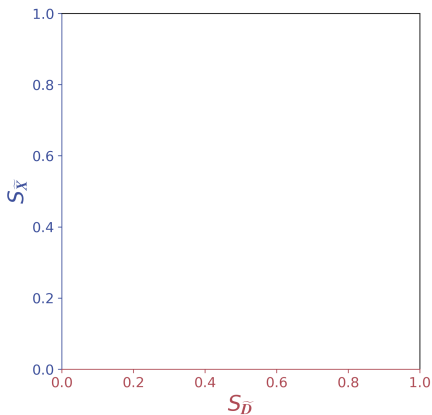


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Interpretation

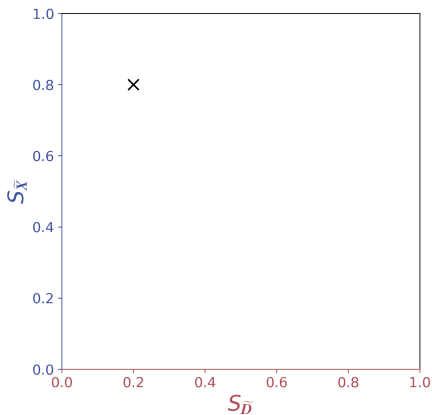


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Interpretation

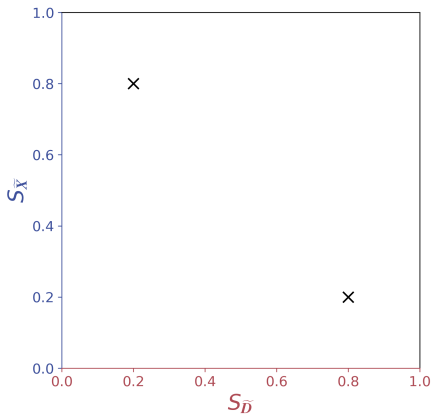


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Interpretation

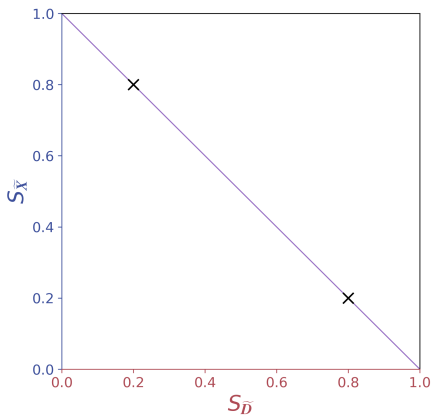


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Interpretation

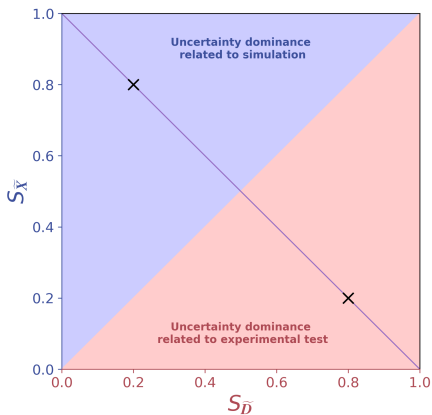


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Interpretation

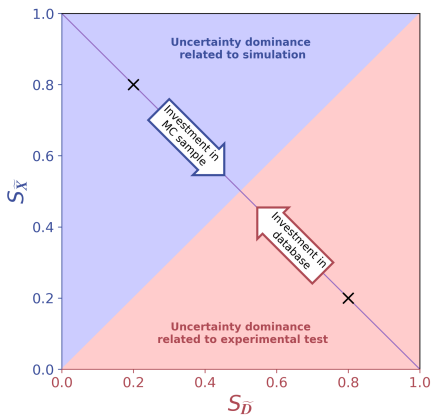


Figure 3.1: Interpretation of Sobol' indices associated to the database and the Monte Carlo sample.

Sensitivity analysis - Independance of inputs

However, the dependance of \tilde{X} to \tilde{D} is an issue for the sensitivity analysis.

Sensitivity analysis - Independance of inputs

However, the dependance of $\tilde{\mathbf{X}}$ to $\tilde{\mathbf{D}}$ is an issue for the sensitivity analysis.

Isoprobabilistic transformation

The transformation $\mathcal{T}_{\mathbf{D}}$ [9, 10, 11] is performed here to work with an independent sample $\tilde{\mathbf{U}} = \{\mathbf{U}^{(j)}, j = 1, \dots, N_{\mathbf{X}}\}$:

$$\mathcal{T}_{\mathbf{D}} : \begin{cases} [0, 1]^d & \longrightarrow \mathcal{X} \\ \mathbf{U} & \longmapsto \mathbf{X} \end{cases}, \quad (3.2)$$

with $\mathbf{U}^{(j)} \overset{i.i.d.}{\sim} \mathcal{U}[0, 1]^d$.

[9] M Rosenblatt. Remarks on a multivariate transformation. The annals of mathematical statistics, 1952.

[10] AE Brockwell. Universal residuals: A multivariate transformation. Statistics probability letters, 2007.

[11] R Lebrun and others. Do rosenblatt and nataf isoprobabilistic transformations really differ? Probabilistic Engineering Mechanics, 2009.

Sensitivity analysis - Solution

Solution B

The investment of data is guided by the highest index:

- $S_{\tilde{D}} > S_{\tilde{U}}$ → Investment in the **database (experimental tests)**,
- $S_{\tilde{D}} < S_{\tilde{U}}$ → Investment in the **MC sample (simulation)**.

Table of contents

- 1 Context
- 2 Proposed approach
- 3 Sensitivity analysis
- 4 Illustration**
- 5 Conclusion

Illustration - Cantilever Beam

Mean deflection of the free end of a cantilever beam:

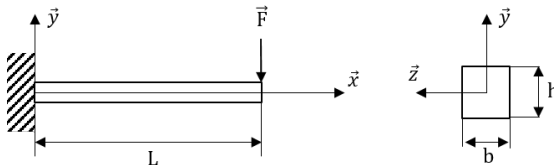


Figure 4.1: Representation of a cantilever beam where F is the transverse load applied on the free end of the beam of length L , Young's modulus E and cross-section bh .

Table 4.1: Probabilistic models associated to independent input variables for a cantilever beam toy-case. [12]

$$\phi(F, L, E, b, h) = \frac{4FL^3}{Ebh^3}$$

Input variable	Distribution	Mean	Coefficient of variation
F	LogNormal	556.8 [N]	0.08
L	Normal	4290 [mm]	0.1
E	LogNormal	2.10^5 [MPa]	0.06
b	Normal	62 [mm]	0.1
h	Normal	98.7 [mm]	0.1

[12] L Baoyu and others. Reliability analysis based on a novel density estimation method for structures with correlations, 2017.

Illustration - Results

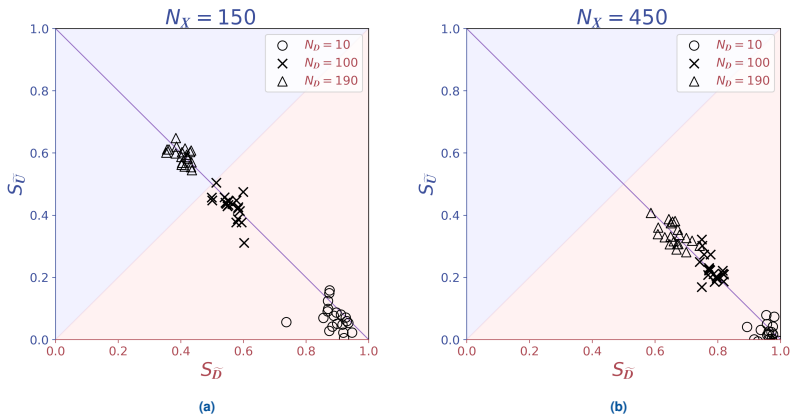


Figure 4.2: Evolution of Sobol' indices for the cantilever beam toy-case at $N_D = [10, 100, 190]$ and (a) $N_X = 150$ (b) $N_X = 450$. Estimation of $n = 20$ indices for each combination.

Table of contents

- 1 Context
- 2 Proposed approach
- 3 Sensitivity analysis
- 4 Illustration
- 5 Conclusion

Conclusion

Framework

- The probabilistic model is **unknown** and is inferred from **experimental tests**,
- A **small-data context** is imposed by costly **experimental tests** and a costly **black box function**.

Conclusion

Framework

- The probabilistic model is **unknown** and is inferred from **experimental tests**,
- A **small-data context** is imposed by costly **experimental tests** and a costly **black box function**.

Current method

- A) Takes into account the **database** uncertainty in the variance of the estimator,
- B) Answers the **test-simulation** trade-off by guiding the investment of data in the driving source of uncertainty.

Conclusion

Framework

- The probabilistic model is **unknown** and is inferred from **experimental tests**,
- A **small-data context** is imposed by costly **experimental tests** and a costly **black box function**.

Current method

- A) Takes into account the **database** uncertainty in the variance of the estimator,
- B) Answers the **test-simulation** trade-off by guiding the investment of data in the driving source of uncertainty.

Perspectives

- Reduction of the computational burden with importance sampling methods, [13]
- Quantification of the amount of data to invest while considering cost differences.

[13] A Owen, Y Zhou. Safe and effective importance sampling. Journal of the American Statistical Association, 2000.

References I

- [1] Gabriel Sarazin.
Analyse de sensibilité fiabiliste en présence d'incertitudes épistémiques introduites par les données d'apprentissage.
PhD thesis, Toulouse, ISAE, 2021.
- [2] James K Lindsey et al.
Parametric statistical inference.
Oxford University Press, 1996.
- [3] Alan Julian Izenman.
Review papers: Recent developments in nonparametric density estimation.
Journal of the american statistical association, 86(413):205–224, 1991.

References II

- [4] Vincent Chabridon.
Analyse de sensibilité fiabiliste avec prise en compte d'incertitudes sur le modèle probabiliste-Application aux systèmes aérospatiaux.
PhD thesis, Université Clermont Auvergne(2017-2020), 2018.
- [5] Chong Ho Yu.
Resampling methods: concepts, applications, and justification.
Practical Assessment, Research, and Evaluation, 8(1):19, 2002.
- [6] Bradley Efron.
The jackknife, the bootstrap and other resampling plans.
SIAM, 1982.

References III

- [7] Ilya M Sobol'.
Sensitivity analysis for non-linear mathematical models.
Mathematical modelling and computational experiment,
1:407–414, 1993.
- [8] Fabrice Gamboa, Alexandre Janon, Thierry Klein, A Lagnoux,
and Clémentine Prieur.
Statistical inference for sobol pick-freeze monte carlo method.
Statistics, 50(4):881–902, 2016.
- [9] Murray Rosenblatt.
Remarks on a multivariate transformation.
The annals of mathematical statistics, 23(3):470–472, 1952.
- [10] Anthony Brockwell.
Universal residuals: A multivariate transformation.
Statistics & probability letters, 77(14):1473–1478, 2007.

References IV

- [11] Régis Lebrun and Anne Dutfoy.
Do rosenblatt and nataf isoprobabilistic transformations really differ?
Probabilistic Engineering Mechanics, 24(4):577–584, 2009.
- [12] LI Baoyu, Leigang Zhang, ZHU Xuejun, YU Xiongqing, and MA Xiaodong.
Reliability analysis based on a novel density estimation method for structures with correlations.
Chinese Journal of Aeronautics, 30(3):1021–1030, 2017.
- [13] Art Owen and Yi Zhou.
Safe and effective importance sampling.
Journal of the American Statistical Association, 95(449):135–143, 2000.