



**HAL**  
open science

# Developments and applications of Shapley effects to reliability-oriented sensitivity analysis with correlated inputs

Marouane Il Idrissi, Vincent Chabridon, Bertrand Iooss

► **To cite this version:**

Marouane Il Idrissi, Vincent Chabridon, Bertrand Iooss. Developments and applications of Shapley effects to reliability-oriented sensitivity analysis with correlated inputs. 2021. hal-03106452v1

**HAL Id: hal-03106452**

**<https://hal.science/hal-03106452v1>**

Preprint submitted on 20 Jan 2021 (v1), last revised 18 May 2021 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Developments and applications of Shapley effects to reliability-oriented sensitivity analysis with correlated inputs

Marouane Il Idrissi<sup>a</sup>, Vincent Chabridon<sup>a,b</sup>, Bertrand Iooss<sup>a,b,c,d</sup>

<sup>a</sup>*EDF Lab Chatou, 6 Quai Watier, 78401 Chatou, France*

<sup>b</sup>*SINCLAIR AI Lab., Saclay, France*

<sup>c</sup>*Institut de Mathématiques de Toulouse, 31062 Toulouse, France*

<sup>d</sup>*Corresponding Author - Email: bertrand.iooss@edf.fr - Phone: +33130877969*

---

## Abstract

Reliability-oriented sensitivity analysis methods have been developed for understanding the influence of model inputs relatively to events characterizing the failure of a system (e.g., a threshold exceedance of the model output). In this field, the target sensitivity analysis focuses primarily on capturing the influence of the inputs on the occurrence of such a critical event. This paper proposes new target sensitivity indices, based on the Shapley values and called “target Shapley effects”, allowing for interpretable influence measures of each input in the case of dependence between the inputs. Two algorithms (a Monte Carlo sampling one, and a given-data algorithm) are proposed for the estimation of these target Shapley effects based on the  $\ell^2$  norm. Additionally, the behavior of these target Shapley effects are theoretically and empirically studied through various toy-cases. Finally, applications on two realistic use-cases (a river flood model and a COVID-19 epidemiological model) are discussed.

*Keywords:* sensitivity analysis, reliability analysis, Sobol’ indices, Shapley effects, input correlation

---

## 1. Introduction

Nowadays, numerical models are intensively used in all industrial and scientific disciplines to describe physical phenomena (e.g., systems of ordinary differential equations in ecosystem modeling, finite element models in structural mechanics, finite volume schemes in computational fluid dynamics) in order to design, analyze or optimize various processes and systems. In addition to this tremendous growth in computational modeling and simulation, the identification and treatment of the multiple sources of uncertainties has become an essential task from the early design stage to the whole system life cycle. As an example, such a task is crucial in the management of complex systems such as those encountered in energy exploration and production [1] and in sustainable resource development [2].

In addition, the emergence of global sensitivity analysis (GSA) of model outputs played a fundamental role in the development and enhancement of these numerical models (see, e.g., [3, 4] for recent reviews). Mathematically, if the model inputs (resp. output) are denoted by  $X$  (resp.  $Y$ ) and the model is written  $G(\cdot)$ , such as

$$Y = G(X), \tag{1}$$

GSA aims at understanding the behavior of  $Y$  with respect to (w.r.t.)  $X = (X_1, \dots, X_d)^\top$  the vector of  $d$  inputs. GSA has been intensively used as a versatile tool to achieve various goals: for instance,

quantifying the relative importance of inputs regarding their influence on the output (a.k.a. "ranking"), identifying the most influential inputs among a large number of inputs (a.k.a. screening) or analyzing the input-output code behavior [5, 6].

When the complex systems are critical or need to be highly safe, numerical models can also be of great help for risk and reliability assessment [7]. Indeed, to track potential failures of a system (which could lead to dramatic environmental, human or financial consequences), numerical models allow to simulate its behavior far from its nominal one (see, e.g., [8] in flood hazard assessment). Analytical or experimental approaches are often out of the question here. Based on these simulations, the tail behavior of the output distribution can be studied and typical *risk measures* can be estimated [9]. Among others, the probability that the output  $Y$  exceeds a given threshold value  $t \in \mathbb{R}$ , given by  $\mathbb{P}(Y > t)$  and often called a *failure probability*, is widely used in many applications. When  $\{Y > t\}$  is a rare event (i.e., associated to a very low failure probability), advanced sampling-based or approximation-based techniques [10] are required to estimate properly the failure probability. In this very specific context, dedicated sensitivity analysis methods have been developed, especially in the structural reliability community (see, e.g., [11, 12, 13]). In such a framework, called *reliability-oriented sensitivity analysis* (ROSA) [14, 15], the idea is to provide importance measures dedicated to the problem of rare event estimation.

To make it more formal, standard GSA methods mostly focus on quantities of interest (QoI) characterizing the central part of the output distribution (e.g., the variance for Sobol' indices [16], the entire distribution for moment-independent indices [17]), while ROSA methods focus on risk measures and their associated practical difficulties (e.g., costly to estimate, inducing a conditioning on the distributions, non-trivial interpretation of the indices). Following [18], ROSA methods can be analyzed regarding the type of study they consider, i.e., according the following two categories:

- *target sensitivity analysis* (TSA) aims at catching the influence of the inputs (considering their entire input domain) on the *occurrence* of the failure event. Basically, this implies to consider the random variable defined by the indicator function  $\mathbb{1}_{\{G(X) > t\}}$  of the failure domain ;
- *conditional sensitivity analysis* aims at studying the influence of the inputs on the *conditional* distribution of the output  $Y|\{G(X) > t\}$ , i.e., exclusively within the critical domain. By Eq. (1), a conditioning also appears on the inputs' domain.

Various indices have been proposed to tackle these two types of studies (see, e.g., [19, 13, 15, 20]). The present paper is dedicated to ROSA (under the assumption that the QoI is a failure probability) and focuses on a TSA study. However, a new issue for TSA is addressed here: the possible statistical dependence between the inputs.

Indeed, most of the common GSA methods (and it is similar for the ROSA ones) have been developed under the assumption of independent inputs. As an example, the well-known Sobol' indices [16] which rely on the so-called functional analysis of variance (ANOVA) and Hoeffding decomposition [21], can be directly interpreted as shares of the output variance that are due to each input and combination of inputs (called "interactions") as long as the inputs are independent.

When the inputs are dependent, the inputs' correlations dramatically alter the interpretation of the Sobol' indices. To handle this issue, several approaches have been investigated in the literature.

For instance, [22] proposed to estimate indices for groups of correlated inputs. However, this approach does not allow to quantify the influence of individual inputs. Amongst other similar works, [23, 24] proposed to extend the functional ANOVA decomposition to a more general one (e.g., taking the covariance into account). However, the indices obtained for these approaches can be negative, which limits their practical use due to interpretability issues. Parallel to this, other works (see, e.g., [25, 26]) considered a Gram–Schmidt procedure to decorrelate the inputs and proposed to estimate two kinds of contributions for each variable (an uncorrelated one and a correlated one). These works finally resulted in the proposition of a set of four Sobol’ indices (instead of the two standard ones which are the first-order index and total index in the independent case) which enable to fully capture the correlation effects in the GSA [27]. Despite this achievement, this approach remains difficult to implement in practice (see [28] for extensive studies). Finally, the VARS approach [29] (allowing a thorough analysis of the inputs-output relationships) can handle inputs’ correlation but is out of scope of the present work which only focuses on variance-based sensitivity indices, directly computed from the numerical model.

Recently, another research track has been developed by considering another type of indices: the *Shapley effects*. The initial formulation originates from the “Shapley values” developed in the field of Game Theory [30, 31]. The underlying idea is to fairly distribute both gains and costs to multiple players working cooperatively. By analogy with the GSA framework, the inputs can be seen as the players while the overall process can be seen as attributing shares of the output variability to the inputs. Considering the variance of the output in a GSA formulation leads to the so-called “Shapley effects” proposed by [32]. In the same vein, [33, 34, 28] bridge the gap between Sobol’ indices and Shapley effects while illustrating the usefulness of these new indices to handle correlated inputs in the GSA framework.

Thus, the present work tries to extend the use of Shapley effects to the ROSA context. To sum up, the idea is to provide a ROSA index enabling to perform TSA (i.e., capturing the influence of the inputs on a risk measure, typically a failure probability here) under the constraint of dependent inputs. Moreover, this work relies on the use of recent promising results and numerical tools (both in field of TSA [35] and Shapley effects’ estimation [36]).

The outline of this paper is the following. Section 2 is devoted to a pedagogical introduction of the statistical dependence issues for variance-based sensitivity indices, that can be solved by Shapley effects. Section 3 presents a new formulation of TSA based on Shapley effects leading to target Shapley effects, while Section 4 develops two algorithms for their estimation. Section 5 provides illustrations on simple toy-cases which give analytical expressions of the target Shapley effects, allowing to deeply study their behavior. Section 6 applies these new sensitivity indices to two use-cases: a simplified model of a river flood and an epidemiological model related to the Covid-19 disease. Finally, Section 7 gives conclusions and research perspectives.

All along the paper, the mathematical notation  $\mathbb{E}(\cdot)$  (resp.  $\mathbb{V}(\cdot)$ ) will represent the expectation (resp. variance) operator.

## 2. Variance-based sensitivity analysis with dependent inputs: the Shapley solution

While devoted to computer experiments, GSA has closed connections with multivariate data analysis and statistical learning [37, 38]. Indeed, in all these topics, one important issue is often to provide a weight to some variables (the inputs) w.r.t. its impact on another variables (the outputs). Depending on the domain, such a weight can be either called a “sensitivity index” or an “importance measure”. A very convenient way is to base these weights on the ANOVA (analysis of variance) decomposition [37, 16] of the output variance. Indeed, such a decomposition provides a natural sharing of the output variance in shares due to each input. The principle of the “variance-based sensitivity indices” [5] consists then in understanding how to separate the contribution of each  $X_i$  in the variance of  $Y$ . However, due to potential statistical dependencies between inputs, this sharing cannot be directly performed. Starting from a simple case such as the linear model, chosen for pedagogical purposes, this section provides a reminder on this topic while illustrating the great interest of Shapley effects in practice.

### 2.1. Understanding the correlation issues via the linear model case

In this section, the aim is to quantify the relative importance of  $d$  scalar inputs  $X_j$  ( $j = 1, \dots, d$ ) by fitting on a data sample (coming from the model Eq. (1)) a linear regression model so as to predict a scalar output  $Y$ :

$$Y(X) = \sum_{j=0}^d \beta_j X_j + \epsilon, \quad (2)$$

where  $X_0 = 1$ ,  $\beta = (\beta_0, \dots, \beta_d)^\top \in \mathbb{R}^{d+1}$  is the effects vector and  $\epsilon \in \mathbb{R}$  the model’s error of variance  $\sigma^2$ . If a sample of inputs and outputs  $(\mathbf{X}^n, \mathbf{Y}^n) = (X_1^{(i)}, \dots, X_d^{(i)}, Y^{(i)})_{i=1, \dots, n}$  is available (with  $n > d$ ), the Ordinary Least Squares method (see, e.g., [37]) can easily be used to estimate the parameters  $\beta$  and  $\sigma^2$  in the linear regression model in Eq. (2). Moreover, one obtains the predictor  $\hat{Y}(x^*)$  of  $Y$  at any prediction point  $x^*$ . An important validation metric of this model is the classical *coefficient of determination* given by:

$$R_{Y(X)}^2 = \sum_{i=1}^n \left[ \hat{Y}(X^{(i)}) - \bar{Y} \right]^2 / \left[ Y^{(i)} - \bar{Y} \right]^2 \quad (3)$$

where  $\bar{Y}$  is the output empirical mean.  $R_{Y(X)}^2$  represents the percentage of output variability explained by the linear regression model of  $Y$  on  $X$ . Finally, from Eq. (2), the variance decomposition expresses as:

$$\mathbb{V}(Y) = \sum_{j=1}^d \beta_j^2 \mathbb{V}(X_j) + 2 \sum_{k>j} \beta_j \beta_k \text{Cov}(X_j, X_k) + \sigma^2. \quad (4)$$

In the specific case of independent inputs, the covariance terms cancel and the standard ANOVA (i.e.,  $\mathbb{V}(Y) = \sum \beta_j^2 \mathbb{V}(X_j) + \sigma^2$ ) is obtained. Then, global sensitivity indices, called Standardized Regression Coefficients (SRC), can be directly computed:

$$\text{SRC}_j = \beta_j \sqrt{\mathbb{V}(X_j) / \mathbb{V}(Y)}. \quad (5)$$

The estimation of the SRC is made by replacing the terms in Eq. (5) by their estimates. Interestingly, this metric for relative importance is signed (thanks to the regression coefficient sign), giving the sense

of variation of the output w.r.t. each input. Moreover,  $\text{SRC}_j^2$  represents a share of variance and the sum of all the  $\text{SRC}_j^2$  approaches  $R^2$  (i.e., the amount of explained variance by the linear model). Note that, in a perfect linear regression model (i.e., without any random error term  $\epsilon$ ),  $\text{SRC}_j$  is equal to the linear Pearson's correlation coefficient between  $X_j$  and  $Y$  (denoted by  $\rho(X_j, Y)$ ). Note also that the ANOVA and  $\text{SRC}^2$  extend to the functional ANOVA and Sobol' indices in the general (non-linear model) case (see Appendix A).

When the inputs are dependent, the main concern is to allocate the covariance terms in Eq. (4) to the various inputs. In this case, the Partial Correlation Coefficient (PCC) has been promoted in GSA [39, 5] as a substitute to the SRC, in order to cancel the effects of other inputs when allocating the weight of one input  $X_j$  in the variance of  $Y$ :

$$\text{PCC}_j = \rho(X_j - \widehat{X}_{-j}, Y - \widehat{Y}_{-j}) \quad (6)$$

where  $X_{-j}$  is the vector of all the  $d$  inputs except  $X_j$ ,  $\widehat{X}_{-j}$  is the prediction of the linear model expressing  $X_j$  w.r.t.  $X_{-j}$  and  $\widehat{Y}_{-j}$  is the prediction of the linear model  $Y$  w.r.t.  $X_{-j}$ . However, PCC is not a right sensitivity index of the input. Indeed, it consists in measuring the linear correlation between  $Y$  and  $X_j$  by fixing  $X_{-j}$ , and is then a measure of the linearity (and not the importance) between the output and one input.

Instead of controlling other inputs  $X_{-j}$  such as done in the PCC, the Semi-Partial Correlation Coefficient (SPCC) quantifies the proportion of the output variance explained by  $X_j$  after removing the information brought by  $X_{-j}$  (on  $X_j$ ) [40]:

$$\text{SPCC}_j = \rho(X_j - \widehat{X}_{-j}, Y) . \quad (7)$$

SPCC can also be expressed by using the relation  $\text{SPCC}_j^2 = R_{Y(X)}^2 - R_{Y(X_{-j})}^2$ , which clearly shows that SPCC gives the additional explanatory power of the input  $X_j$  in the linear regression model of  $Y$  on  $X$ . However, the SPCC of highly correlated inputs will be small, despite their "real" explanatory power on the output. This aspect seems to be the main drawback of SPCC and probably explains its lack of popularity for GSA purposes.

In order to give an intuitive view of the *multicollinearity* issue (i.e., multiple linear regression with correlated inputs), we use Venn diagrams (see Figure 1), by considering two inputs  $X_1$  and  $X_2$  and one output  $Y$ . From Figure 1, the coefficient of determination can be written as:

$$R_{Y(X_1, X_2)}^2 = \frac{a + b + c}{a + b + c + \sigma^2} , \quad (8)$$

where  $a + b + c + \sigma^2$  is equal to the variance of  $Y$  and  $a + b + c$  represents the part of explained variance by the regression model (with  $b = 0$  in the uncorrelated case). In this elementary example, the previously introduced sensitivity indices are given by [41]:

$$\begin{aligned} \text{SRC}_1^2 &= (a + b)/(a + b + c + \sigma^2) , & \text{SRC}_2^2 &= (c + b)/(a + b + c + \sigma^2) , \\ \text{PCC}_1^2 &= a/(a + \sigma^2) , & \text{PCC}_2^2 &= c/(c + \sigma^2) , \\ \text{SPCC}_1^2 &= a/(a + b + c + \sigma^2) , & \text{SPCC}_2^2 &= c/(a + b + c + \sigma^2) . \end{aligned} \quad (9)$$

Thus, one can understand the limitations of SRC, PCC and SPCC when correlation is present: the variance share which comes from the correlation between inputs (i.e., the  $b$  value in Figure 1 - right) is allocated two times with the SRC but not allocated at all with SPCC, while PCC does not represent any variance sharing.

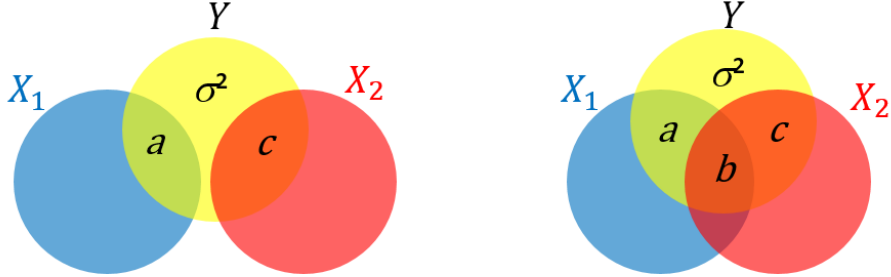


Figure 1: Inspired from [41]. Illustration scheme of the effect of two inputs  $X_1$  and  $X_2$  on an output variable  $Y$  when they are: uncorrelated (left) or correlated (right).

The three problems above can be solved by using another sensitivity index which finds a way to partition the  $R^2$  among the  $d$  inputs: the LMG [42, 43] (acronym based on the authors' names, i.e., "Lindeman - Merenda - Gold") uses sequential sums of squares from the linear model and obtains an overall measure by averaging over all orderings of inputs. Mathematically, let  $u$  be a subset of indices in the set of all subsets of  $\{1, \dots, d\}$  and  $X_u = (X_j : j \in u)$  a group of inputs. LMG is based on the measure of the elementary contribution of any given variable  $X_j$  to a given subset model  $Y(X_u)$  by the increase in  $R^2$  that results from adding that predictive variable to the regression model:

$$\text{LMG}_j = \frac{1}{d!} \sum_{\substack{\pi \in \text{permutations} \\ \text{of } \{1, \dots, d\}}} \left[ R_{Y(X_{v \cup \{j\}}}^2) - R_{Y(X_v)}^2 \right] \quad (10)$$

with  $v$  the indices entered before  $j$  in the order  $\pi$ . In Eq. (10), the sum is performed over all the permutations of  $\{1, \dots, d\}$ . For the case of two inputs (see Figure 1), we can easily show that:

$$\text{LMG}_1 = (a + b/2)/(a + b + c + \sigma^2), \quad \text{LMG}_2 = (c + b/2)/(a + b + c + \sigma^2). \quad (11)$$

Then, in the LMG framework, the  $R_{Y(X_1, X_2)}^2$  has been perfectly shared into two parts with an equitable distribution of the  $b$  term between  $X_1$  and  $X_2$ .

This allocation principle exactly corresponds to the application of the *Shapley values* [30] on the linear model. This attribution method has been primarily used in cooperative game theory, allowing for a cooperative allocation of resources between players based on their collective production (see Appendix B for a more formal definition). The Shapley values solution consists in fairly distributing both gains and costs to several actors working in coalition. In situations when the contributions of each actor are unequal, it ensures that each actor gains as much or more as they would have from acting independently. Now, if the actors are identified with a set of inputs and the value assigned to each coalition is identified to the explanatory power of the subset of model inputs composing the coalition, one obtains the LMG in Eq. (10).

## 2.2. Shapley effects

In the general case with no hypothesis on the form of the model  $G(\cdot)$  (Eq. (1)), variance-based sensitivity indices have been extensively developed [16, 5] and applied for GSA of complex models (see, e.g., [44]). Indeed, in the independent inputs' case, it allows a variance decomposition of the model output in different shares (called ‘‘Sobol’ indices’’) induced by each input and each possible interaction between inputs in the model (see Appendix A for the theoretical details).

For the dependent inputs' case, ideas coming from game theory and Shapley values (as for LMG in Eq. (10)), have been recently introduced [32] in order to measure input influence by

$$Sh_j = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} (\text{val}(A \cup \{j\}) - \text{val}(A)), \quad (12)$$

where  $\text{val}(A)$  is the so-called *value* function (which is, somehow a cost function) assigned to a subset  $A \in \mathcal{P}_d$  of inputs,  $\mathcal{P}_d$  is the set of all possible subsets of  $\{1, \dots, d\}$ ,  $\{-j\}$  denotes the set of indices  $\{1, \dots, d\} \setminus j$  and  $|A|$  is the cardinal of  $A$ .

For GSA purposes, [32] proposes to use the so-called ‘‘closed Sobol’ indices’’ as the value function:

$$\text{val}(A) = S_A^{\text{clos}} = \frac{\mathbb{V}(\mathbb{E}[G(X) | X_A])}{\mathbb{V}(G(X))}, \quad (13)$$

where  $X_A$  denotes the subset of inputs selected by the indices in  $A$  ( $X_A = (X_i)_{i \in A}$ ). The attribution properties of the Shapley values applied to this particular cost function, leads to the definition of the *Shapley effects*:

$$Sh_j = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} (S_{A \cup \{j\}}^{\text{clos}} - S_A^{\text{clos}}) \quad (14)$$

These Shapley effects allow for a quantification of influence for each input, which intrinsically takes into account both interaction and dependence. Moreover, two important properties of the Shapley effects ( $Sh_j, j = 1, \dots, d$ ) allows for an easy interpretation: they sum up to one and are non-negative. They allow for input ranking, in terms of influence, by allocating each input a percentage of the model output's variance. These indices have been extensively studied in [33, 34]. An alternate way of defining the Shapley effects has been proposed, by taking the following cost function:

$$\text{val}(A) = \frac{\mathbb{E}[\mathbb{V}(G(X) | X_{\bar{A}})]}{\mathbb{V}(G(X))} \quad (15)$$

where  $\bar{A} = \{1, \dots, d\} \setminus A$  which lead to the equivalent definition of the Shapley effects in Eq. (14). This results allows for alternate estimation methods, as outlined in [45].

In order to illustrate the allocation system of the Shapley effects, one can first consider a model



with three inputs  $X = (X_1, X_2, X_3)^\top$ . From Eq. (14), one gets:

$$\begin{aligned} Sh_1 &= \frac{1}{3}S_1^{\text{clos}} \\ &+ \frac{1}{6}[(S_{\{1,2\}}^{\text{clos}} - S_2^{\text{clos}}) + (S_{\{1,3\}}^{\text{clos}} - S_3^{\text{clos}})] \\ &+ \frac{1}{3}(S_{\{1,2,3\}}^{\text{clos}} - S_{\{2,3\}}^{\text{clos}}). \end{aligned}$$

In the case where the three inputs are independent, one obtains:

$$Sh_1 = S_1 + \frac{1}{2}S_{\{1,2\}} + \frac{1}{2}S_{\{1,3\}} + \frac{1}{3}S_{\{1,2,3\}}$$

where one can notice that the Shapley values quantification (in the independent case) consists of the initial Sobol' index of the studied input, plus an equal share of the interaction effects between all the involved inputs. However, if dependence between inputs is assumed, this behavior cannot be clearly illustrated, except in the linear case (see Subsection 2.1).

The Shapley increment ( $S_{A \cup \{j\}}^{\text{clos}} - S_A^{\text{clos}}$ ) can be interpreted as being a quantification of the *residual effects of the input  $j$  in relation to the subset of variables  $A$* . If  $S_A^{\text{clos}}$  is believed to contain the initial effects of  $A$ , plus their interaction effects, and any effects due to the dependence of the inputs in  $A$ , then the Shapley increment quantifies the initial effect of the input  $j$ , its interaction effect with  $A$ , and the effects due to its dependence with the inputs in  $A$ . Then, the Shapley attribution system weights all these residual effects, in order to equally redistribute the interaction or dependence effects between the involved inputs, in the same fashion as the LMG in the linear model case, as depicted in Section 2.1.

It is important to note that the above mentioned behavior of the Shapley effects cannot be verified for any type of dependence structure, due to the lack of a univocal functional decomposition in the case of dependent inputs. However, this behavior has been highlighted in analytical cases in [34] (mainly for Gaussian inputs with a linear correlation structure).

### 3. Reliability-oriented Shapley effects for target sensitivity analysis

#### 3.1. A brief overview of reliability-oriented sensitivity analysis

When focusing on complex systems, one often needs to prevent from possible critical events, which have a low probability of occurring but are associated to a failure of the system. Such a failure might have possible dramatic consequences regarding various factors (e.g., human, environmental, financial). Such a task is covered by the fields of reliability assessment and risk analysis [7, 8]. Mathematically, such a problem implies to focus on a *risk measure* computed from the tail of the output distribution [9]. Performing sensitivity analysis in such a context requires to use dedicated tools which have been gathered by various authors under the denomination of “reliability-oriented sensitivity analysis” (ROSA) (see, e.g., [15, 46, 20]). A large panel of ROSA methods have been proposed in the structural reliability community such as, for example, several variance-based approaches (see, e.g., [47, 13, 15, 48]) and moment-independent approaches (see, e.g., [49, 19, 46]). From the GSA community, several extensions have also been proposed as extensions to handle risk or reliability measures. Among others,

one can mention, for instance, the contrast-based indices proposed by [50] as a versatile tool to handle several types of QoI and then studied by [51, 52] for quantile-oriented formulations, the quantile-based global sensitivity measures [53] or, finally, other indices related to dependence measures [18, 20].

In the context of reliability assessment, a typical risk measure is the *failure probability* given by:

$$p_t^Y \stackrel{\text{def}}{=} \mathbb{P}(Y > t) = \mathbb{P}(G(X) > t) = \mathbb{E}[\mathbb{1}_{\{G(X) > t\}}(X)] = \mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X)] \quad (16)$$

where  $t$  represents a certain scalar threshold value characterizing the state of the system. Typically, the event  $\{Y > t\}$  denotes the failure event. As for  $\mathcal{F}_t$ , it represents the input failure domain, i.e.,  $\mathcal{F}_t \stackrel{\text{def}}{=} \{X \mid G(X) > t\}$ . Thus, performing a ROSA study induces a few challenges: firstly, the variable of interest here is no more  $Y$ , but can be, for instance, a binary event whose occurrence is characterized by the indicator function  $\mathbb{1}_{\mathcal{F}_t}(X)$ ; secondly, this event is likely to be a “rare event” associated to a low failure probability which might be difficult to estimate in practice [10]; thirdly, the type of study one desires to perform has to be reinterpreted regarding the new QoI. Regarding this last point, [18] proposes to focus on two types of studies when dealing with critical events: the first one, called *target sensitivity analysis* (TSA), aims at catching the influence of the inputs on the occurrence of the failure event, while the second one, called “conditional sensitivity analysis” aims at studying the influence of the inputs once the threshold value has been reached (i.e., within the failure domain). The present paper is dedicated to ROSA (under the assumption that the QoI is a failure probability given by Eq. (16)) and focuses on a TSA study.

To illustrate this in plain text, one can use the example of the modeling of the water level in a river protected by a dyke. From the traditional GSA point of view, the central question would be “Which inputs influence the water level?”, while in the TSA paradigm, one focuses more on the question “Which inputs influence the occurrence of a specific flood level?”. Note that this particular example is more specifically studied in Subsection 6.1.

In the case of independent inputs, a first category of sensitivity indices dedicated to TSA is the “target Sobol’ indices” whose first formulation has been proposed by [19]. Similarly, one uses here the closed Sobol’ index (see Appendix A) as follows:

$$\text{T-S}_A = \frac{\mathbb{V}(\mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X) \mid X_A])}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))}. \quad (17)$$

where  $\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X)) = p_t^Y(1-p_t^Y)$ . Several estimation schemes for these indices have been proposed in the rare event context [13, 15, 54]. To illustrate the behavior of this type of basic index, one can consider the additive model given by  $Y = X_1 + X_2 + X_3$ , with  $X = (X_1, X_2, X_3)^\top$ , three independent standard Gaussian random variables. The left plot of Figure 2 represents the probability density function (pdf) of  $Y$ , with four different values for the threshold  $t$ , corresponding to four different failure probability levels. The right plot of Figure 2 presents the different values of  $\text{T-S}_A$ . Note that, for this specific example, the second-order indices verify  $\text{T-S}_{\{1,2\}} = \text{T-S}_{\{1,3\}} = \text{T-S}_{\{2,3\}}$ . Moreover, one can remark that the more  $t$  is restrictive or loose (i.e., the failure probability being “close” to 0 or 1), the more the third-order closed Sobol’ index for TSA increases, indicating high interaction effects between the three inputs. Note that the understandings of this behavior falls under the conditional sensitivity analysis paradigm, which is out of the scope of this paper. However, the acknowledgment of this phenomenon

remains important for better understanding the behavior of the new indices proposed further in the paper.

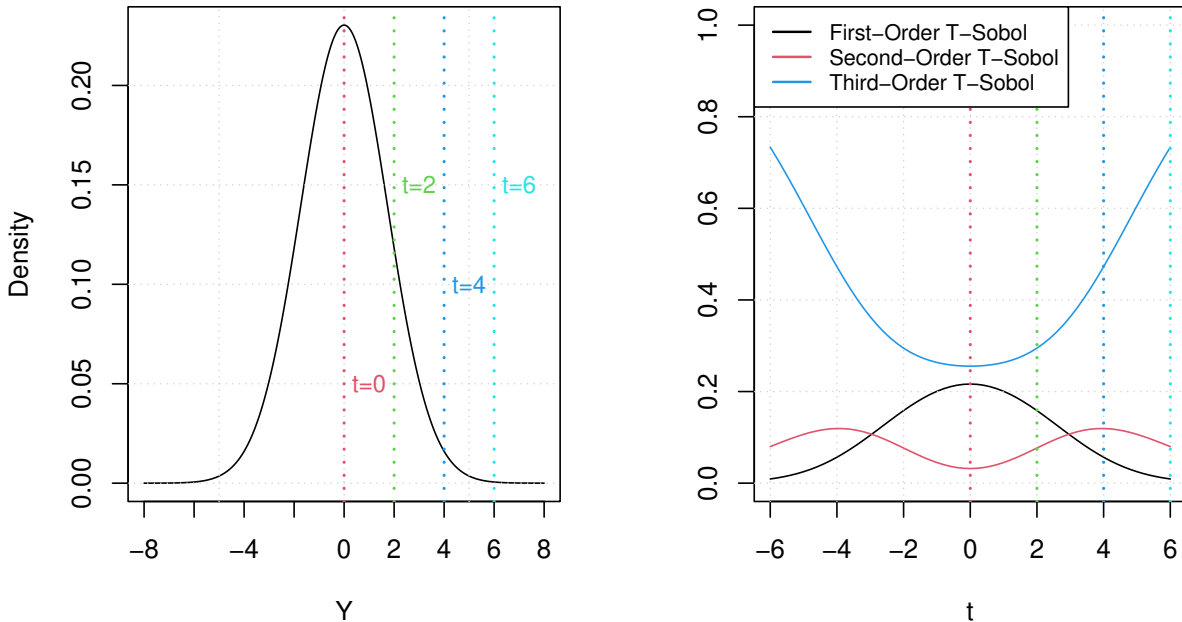


Figure 2: Probability density function of the output with four different threshold values (left) and the related target Sobol' indices (right) for  $Y$  being the sum of three independent Gaussian random variables.

Another category of sensitivity indices dedicated to TSA gathers the “moment-independent” ROSA indices. Among others, one can mention the two indices proposed by [49] which are given by:

$$\eta_A = \frac{1}{2} \mathbb{E} \left[ \left| p_t^Y - p_t^{Y|X_A} \right| \right] \quad (18a)$$

$$\delta_A = \frac{1}{2} \mathbb{E} \left[ \left( p_t^Y - p_t^{Y|X_A} \right)^2 \right]. \quad (18b)$$

where  $p_t^{Y|X_A}$  denotes the conditional failure probability when  $X_A$  is fixed. Note that, if the first index does not require any independence assumption for quantifying input influence, it is known to be difficult to estimate in practice [46]. As for the second one, it is simply proportional to the target Sobol' index given in Eq. (17). Note that an extension of the  $\delta_A$  index has been proposed in [55] for correlated inputs. It relies on a similar orthogonalization procedure strategy as proposed by [26]. However, as mentioned previously, this tends to increase the number of indices to estimate so as to properly interpret the results.

The following section develops the *distance-based TSA indices*, highlighting their links with existing TSA indices and proposing new TSA indices inspired from the Shapley framework presented in Subsection 2.2.

### 3.2. Distance-based TSA indices

As recalled by several authors [50, 18], one can notice that  $\mathbb{V}(\mathbb{E}[Y|X_A]) = \mathbb{E}[(\mathbb{E}[Y|X_A] - \mathbb{E}[Y])^2]$ . This equality can be interpreted as the expected squared distance between two expectations, and thus

allows to apprehend Sobol' indices (see Eq. (13)) as distance-based indices. Such a general point of view has been adopted by [50] to provide a generalization of the Sobol' indices using *contrast functions*.

By applying a similar idea in the TSA paradigm, one can extend the standard  $T-S_A$  and  $\eta_A$  indices to more general cases using statistical distances. Thus, one can define the distance-based TSA index, relative to a subset of inputs  $A \in \mathcal{P}_d$ , as follows:

$$T-S_A^{\mathcal{D}} = \frac{\mathbb{E}\left[\mathcal{D}\left(p_t^Y, p_t^{Y|X_A}\right)\right]}{\mathbb{E}\left[\mathcal{D}\left(p_t^Y, p_t^{Y|X}\right)\right]} \quad (19)$$

where  $\mathcal{D}(\cdot, \cdot)$  denotes any distance function. Links can be made between these indices and the one previously presented, through the use of specific distance functions. For example, by choosing the distance derived from the  $\ell^1$  norm (i.e., the absolute difference), one can find that the corresponding distance-based TSA index is proportional to the  $\eta_A$  index, that is:

$$T-S_A^{\ell^1} = \frac{\mathbb{E}\left[|\mathbb{E}[\mathbb{1}_{F_t}(X)] - \mathbb{E}[\mathbb{1}_{F_t}(X)|X_A]|\right]}{\mathbb{E}\left[|\mathbb{1}_{\mathcal{F}_t}(X) - \mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X)]|\right]} = \frac{2}{\mathbb{E}\left[|\mathbb{1}_{\mathcal{F}_t}(X) - \mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X)]|\right]} \eta_A \quad (20)$$

Moreover, if using the distance derived from the  $\ell^2$  norm (i.e., the squared difference), one can remark the related distance-based TSA indices are equal to the closed Sobol' index for TSA, as defined in Eq. (17):

$$T-S_A^{\ell^2} = T-S_A = \frac{\mathbb{V}\left(\mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X) | X_A]\right)}{\mathbb{V}\left(\mathbb{1}_{\mathcal{F}_t}(X)\right)}. \quad (21)$$

This framework of distance-based TSA indices allows to quantify the influence of inputs (or subsets of inputs), but the previously introduced indices are not suited for the case of dependent inputs. However, such a framework enables to produce relevant candidates as cost functions for the Shapley setting as presented in the following.

### 3.3. ( $\mathcal{D}$ )-target Shapley effects

In this subsection, a novel family of TSA indices is proposed, namely the ( $\mathcal{D}$ )-target Shapley effects. As mentioned above, these indices are built by taking the distance-based TSA indices, defined in Eq. (19), as cost functions in a Shapley attribution procedure (see Eq. (12)). For a specific input  $j \in \{1, \dots, d\}$ , its ( $\mathcal{D}$ )-target Shapley effects can be defined as being:

$$T-Sh_j^{\mathcal{D}} = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} \left( T-S_{A \cup \{j\}}^{\mathcal{D}} - T-S_A^{\mathcal{D}} \right) \quad (22)$$

where  $\{-j\} = \{1, \dots, d\} \setminus j$ . The main property allowing for a clear interpretation of the ( $\mathcal{D}$ )-target Shapley effects is the following:

**Property 1** (( $\mathcal{D}$ )-target Shapley effects decomposition). *Let  $A \in \mathcal{P}_d$ , and  $val(A) = T-S_A^{\mathcal{D}}$ . For any*

distance function  $\mathcal{D}(\cdot, \cdot)$ , the following property holds:

$$\sum_{j=1}^d T\text{-Sh}_j^{\mathcal{D}} = 1. \quad (23)$$

It is important to note that this decomposition property does not rely on any independence assumption on the inputs of the numerical model. However, in order to ensure a meaningful interpretation of these indices, (e.g., such as a percentage of a meaningful statistic in terms of input importance), one needs to ensure that the  $T\text{-Sh}_j$  are non-negative, for all  $j = 1, \dots, d$ .

By choosing the cost function being equal to  $T\text{-S}_A^{\ell^1}$  (i.e.,  $\mathcal{D}(x, y) = |x - y|$ ), one can define the  $(\ell^1)$ -target Shapley effect associated to a variable  $j \in \{1, \dots, d\}$  as being:

$$T\text{-Sh}_j^{\ell^1} = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} \left( T\text{-S}_{A \cup \{j\}}^{\ell^1} - T\text{-S}_A^{\ell^1} \right). \quad (24)$$

These indices are non-negative (see the proof in Appendix C.1) which allows the  $(\ell^1)$ -target Shapley effects to be interpreted as the percentage of the mean absolute deviation of the indicator function (i.e.,  $\mathbb{E} \left[ \left| \mathbb{1}_{\mathcal{F}_t}(X) - \mathbb{E}[\mathbb{1}_{\mathcal{F}_t}(X)] \right| \right]$ ) one can allocate to each input  $X_j$  with  $j \in \{1, \dots, d\}$ .

By choosing  $T\text{-S}_A^{\ell^2}$  as the cost function in the Shapley framework (i.e.,  $\mathcal{D}(x, y) = (x - y)^2$ ), the  $(\ell^2)$ -target Shapley effect associated to the variable  $j \in \{1, \dots, d\}$  can be defined as:

$$T\text{-Sh}_j^{\ell^2} = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} \left( T\text{-S}_{A \cup \{j\}}^{\ell^2} - T\text{-S}_A^{\ell^2} \right) \quad (25)$$

Being also non-negative (see the proof in Appendix C.2), they can be interpreted as a percentage of the variance of the indicator function allocated to the input  $X_j$  with  $j \in \{1, \dots, d\}$ . Moreover, using Eq. (21), by analogy with Eqs. (13) and (15), if one chooses to define the cost function as being:

$$\text{val}(A) = T\text{-E}_A \stackrel{\text{def}}{=} \frac{\mathbb{E} \left[ \mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X) | X_{\bar{A}}) \right]}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))} \quad (26)$$

with  $\bar{A} = \{1, \dots, d\} \setminus A$ , then one has an equivalent way of defining the  $(\ell^2)$ -target Shapley effect.

In the following, by analogy to  $T\text{-S}_A^{\ell^2} = T\text{-S}_A$  and to simplify the terminology and notations, the  $(\ell^2)$ -target Shapley effect  $T\text{-Sh}_j^{\ell^2}$  will be called the target Shapley effect and denoted  $T\text{-Sh}_j$ :

$$T\text{-Sh}_j \stackrel{\text{def}}{=} T\text{-Sh}_j^{\ell^2}. \quad (27)$$

#### 4. Estimation methods and practical implementation of target Shapley effects

The estimation of the target Shapley effects Eq. (25) can be split into two steps:

- **Step #1:** estimation of the *conditional elements*, i.e., the estimation of  $T\text{-S}_A$  or  $T\text{-E}_A$  for all  $A \in \mathcal{P}_d$ ;

- **Step #2:** an *aggregation procedure*, i.e., a step to compute the T-Sh<sub>j</sub> by plugging in the previous estimations of Step #1 in Eq. (25).

In the following, two estimation methods are proposed: the first one based on a Monte Carlo sampling procedure, and the second one based on a nearest-neighbor approximation technique.

#### 4.1. Monte Carlo sampling-based estimation

This procedure, introduced in [45] for the estimation of Shapley effects, relies on a Monte Carlo estimation of the conditional elements. It requires the knowledge and the ability to sample from the marginal distributions of the inputs, that is  $P_{X_A}$  for all  $A \subseteq \{1, \dots, d\} \setminus \emptyset$ , as well as from all the conditional distributions, that is  $P_{X_{\bar{A}}|X_A}$ , for all possible subsets of inputs  $A$ . Additionally, one also needs to be able to evaluate the model  $G(\cdot)$  which is mostly the case in the context of uncertainty quantification of numerical computer models (ignoring the potential difficulties related to the cost of a single evaluation of  $G(\cdot)$ ) [1].

In order to estimate a conditional element T-S<sub>A</sub>, one needs to draw several i.i.d. samples:

- an i.i.d. sample of size  $N$  drawn from  $P_X$  and denoted by  $(X^{(1)}, \dots, X^{(N)})$ ;
- another i.i.d. sample of size  $N_v$  drawn from  $P_{X_A}$  and denoted by  $(X_A^{(1)}, \dots, X_A^{(N_v)})$ ;
- for each element  $X_A^{(i)}, i = 1, \dots, N_v$ , a corresponding sample of size  $N_p$  drawn from  $P_{X_{\bar{A}}|X_A}$  given that  $X_A = X_A^{(i)}$  and denoted by  $(\tilde{X}_i^{(1)}, \dots, \tilde{X}_i^{(N_p)})$ .

Then, the Monte Carlo estimator of T-S<sub>A</sub> can be defined as:

$$\widehat{\text{T-S}}_{A,\text{MC}} = \frac{\sum_{i=1}^{N_v} \left( \frac{1}{N_p} \sum_{j=1}^{N_p} \mathbb{1}_{\mathcal{F}_t}(\tilde{X}_i^{(j)}, X_A^{(i)}) - \hat{p}_t^Y \right)^2}{(N_v - 1) \hat{p}_t^Y (1 - \hat{p}_t^Y)} \quad (28)$$

with

$$\hat{p}_t^Y = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\mathcal{F}_t}(X^{(i)}). \quad (29)$$

Finally, the aggregation procedure gives:

$$\widehat{\text{T-Sh}}_{j,\text{MC}} = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} \left( \widehat{\text{T-S}}_{A \cup \{j\},\text{MC}} - \widehat{\text{T-S}}_{A,\text{MC}} \right). \quad (30)$$

Thus, one gets that  $\widehat{\text{T-Sh}}_{j,\text{MC}}$  is an unbiased consistent estimator of T-Sh<sub>j</sub>.

Algorithm 1 provides a detailed description on how to implement this estimator in practice. This estimation method requires  $(N + d! \times (d - 1) \times N_v \times N_p)$  calls to the numerical model  $G(\cdot)$ . As expected, this first estimation method can become quite expensive in practice. Moreover, numerical models usually encountered in industrial studies can be costly-to-evaluate, which can become a strong limitation for the use of this method.

Another algorithm has been proposed in [45], by leveraging an equivalent definition of the Shapley attribution system, as an arithmetic mean over all the  $d!$  permutations of  $\{1, \dots, d\}$ . In the same

---

**Algorithm 1:** Target Shapley effects estimation by a Monte Carlo procedure.

---

**Input:**  $G, t, d, N, N_v, N_p, \text{simJoint}, \text{simMarginal}, \text{simConditional}$

**Output:**  $(\widehat{\text{T-Sh}}_{j,\text{MC}})_{j=1,\dots,d}$

```

1  $(X^{(1)}, \dots, X^{(N)}) \leftarrow \text{sim\_joint}(N)$  /* Sample from the joint distribution */
2  $\widehat{p}_t^Y \leftarrow \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{G(X^{(i)}) > t\}}(X^{(i)})$  /* Estimate the failure probability */
3 for  $A \in \mathcal{P}_d$  do /* For every subsets of inputs */
4      $(X_A^{(1)}, \dots, X_A^{(N_v)}) \leftarrow \text{simMarginal}(A, N_v)$  /* Sample from the marginal distribution */
5     for  $i = 1, \dots, N_v$  do /* For every element of the marginal distribution sample */
6          $(\widetilde{X}_i^{(1)}, \dots, \widetilde{X}_i^{(N_p)}) \leftarrow \text{simConditional}(\overline{A}, N_p, X_A^{(i)})$  /* Sample from the conditional distribution given the element of the marginal distribution */
7          $\widehat{\text{T-S}}_A \leftarrow \frac{1}{N_v-1} \sum_{i=1}^{N_v} \left( \frac{1}{N_p} \sum_{j=1}^{N_p} \mathbb{1}_{\mathcal{F}_t}(\widetilde{X}_i^{(j)}, X_A^{(i)}) - \widehat{p}_t^Y \right)^2 \times \frac{1}{\widehat{p}_t^Y(1-\widehat{p}_t^Y)}$  /* Compute the conditional element */
8     for  $j = 1, \dots, d$  do /* Aggregation step */
9          $\widehat{\text{T-Sh}}_{j,\text{MC}} \leftarrow 0$ 
10        for  $A \subset \{-j\}$  do /* Apply the Shapley weights to every computed increments */
11             $\widehat{\text{T-Sh}}_{j,\text{MC}+} = \frac{1}{d} \binom{d-1}{|A|}^{-1} \left( \widehat{\text{T-S}}_{A \cup \{j\}} - \widehat{\text{T-S}}_A \right)$ 

```

---

fashion as in Eq. (10), it writes:

$$\widehat{\text{T-Sh}}_j = \frac{1}{m} \sum_{\substack{\pi \in \text{permutations} \\ \text{of } \{1, \dots, d\}}} \left( \widehat{\text{T-S}}_{v \cup \{j\}, \text{MC}} - \widehat{\text{T-S}}_{v, \text{MC}} \right) \quad (31)$$

with  $v$  being the indices before  $j$  in the order  $\pi$ . In Eq. (31), the sum is not performed over all the permutations of  $\{1, \dots, d\}$  but only on  $m$  randomly chosen permutations. By sampling  $m < d!$  permutations, one can drive the computational cost of this algorithm to  $(N + m \times (d-1) \times N_v \times N_p)$  calls to  $G(\cdot)$ , for a less precise, but still convergent estimator.

#### 4.2. Given-data estimation

A “given-data” estimation method has been introduced in [36] to get the Shapley effects. This method can be seen as an extension of the Monte Carlo estimator when only a single i.i.d. input-output sample is available. This method is appropriate when the input distributions are not known or when the numerical model  $G(\cdot)$  is no more available. The main idea behind this method is to replace the exact samples from the conditional distributions  $P_{X_{\overline{A}}|X_A}$  by approximated ones based on a non-parametric nearest-neighbor procedure.

Let  $(X^{(1)}, \dots, X^{(N)})$  be an i.i.d. sample of the inputs  $X$  and  $A \in \mathcal{P}_d \setminus \{\emptyset, [1 : d]\}$ . Let  $k_N^A(l, n)$

be the index such that  $X_A^{(k_N^A(l,n))}$  is the  $n$ -th closest element to  $X_A^{(l)}$  in  $(X_A^{(1)}, \dots, X_A^{(N)})$ . Note that, if two observations are at an equal distance from  $X_A^{(l)}$ , then one of the two is uniformly randomly selected. Finally, one can define an estimator of the equivalent cost function defined in Eq. (26):

$$\widehat{\text{T-E}}_{A,\text{KNN}} = \frac{\sum_{l=1}^N \left( \frac{1}{N_s-1} \sum_{i=1}^{N_s} \left[ \mathbb{1}_{\mathcal{F}_t} \left( X^{(k_N^A(l,i))} \right) - \frac{1}{N_s} \sum_{h=1}^{N_s} \mathbb{1}_{\mathcal{F}_t} \left( X^{(k_N^A(l,h))} \right) \right]^2 \right)}{N \widehat{p}_t^Y (1 - \widehat{p}_t^Y)}. \quad (32)$$

Under some mild assumptions, [36] showed that this estimator does asymptotically converge towards  $\text{T-E}_A$ . With estimates for the conditional elements, one can then define the following plug-in estimator:

$$\widehat{\text{T-Sh}}_{j,\text{KNN}}^1 = \frac{1}{d} \sum_{A \subset \{-j\}} \binom{d-1}{|A|}^{-1} \left( \widehat{\text{T-E}}_{A \cup \{j\}, \text{KNN}} - \widehat{\text{T-E}}_{A, \text{KNN}} \right) \quad (33)$$

where  $\widehat{p}_t^Y$  is the empirical mean of  $\mathbb{1}_{\mathcal{F}_t}(X)$  on the i.i.d. sample. Algorithm 2 represents the procedure for this given-data estimator. This method is less computationally expensive (in terms of model evaluations) compared to the Monte Carlo sampling-based method, since no additional model evaluation, other than the ones in the i.i.d. sample, is required in order to produce estimates of the target Shapley effects. Since the samples of the conditional and marginal distributions are approximated by a non-parametric procedure, this method also allows to reduce the possible input modeling error (e.g., in the context of ill-defined input distributions), at the cost of less accurate estimates. Another constraint is due to the fact that the input-output sample has to be i.i.d. which prevents from its use with, for instance, advanced design of computer experiments.

In [36], a random permutation algorithm, homologous to Eq. (31), has been developed, which allows for reducing the overall complexity of the method, which, for the sake of conciseness, is not developed in this paper.

#### 4.3. Software and reproducibility of results

The algorithms described in the preceding subsections have been implemented in the `sensitivity` R package [56]. More precisely, the `shapleyPermEx()` (sampling-based algorithm) and `sobolshap_knn()` (given-data algorithm) functions can be directly used for the estimation of the target Shapley effects. In the applications of Section 6, only the `sobolshap_knn()` function will be used for numerical tractability. Appendix D provides some minimal code examples for the implementation of the aforementioned algorithms, along with their random permutation variants [36].

All further results can be accessed on a GitLab platform<sup>1</sup>, along with the data used in the following sections. R code files are available, with explicit code, along with all custom-made functions, in order to reproduce the analyses presented in this paper. The procedures for the theoretical approximations of Section 5 are made available, along with the data-simulation functions for the flood case in Subsection 6.1. The two datasets used for Subsection 6.2 are also available. Finally, all the figures can be reproduced by simply re-running the different `RMarkdown` files in the aforementioned GitLab repository.

<sup>1</sup>[https://gitlab.com/milidris/review\\_l2tse](https://gitlab.com/milidris/review_l2tse)



---

**Algorithm 2:** Target Shapley effects estimation by a nearest-neighbor procedure.

---

```

Input:  $X, Y, t$ 
Output:  $(\widehat{\text{T-Sh}}_{j,\text{KNN}})_{j=1,\dots,d}$ 
/* Estimate the failure probability */
1  $\widehat{p}_t^Y \leftarrow \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{G(X^{(i)}) > t\}}(X^{(i)})$  /*
/* For every subsets of inputs */
2 for  $A \in \mathcal{P}_d$  do /*
/* Sample of  $X_A$  */
3  $X_A \leftarrow (X_i^{(j)})_{j=1,\dots,n}^{i \in A}$  /*
4 for  $i = 1, \dots, N$  do /*
/* For each row  $i$  of  $X_A$ , find the  $N_s$  nearest rows in  $X$  */
5  $(\widetilde{X}_i^{A,(j)})_{j=1,\dots,N_s} \leftarrow \text{KNN}(X_A^{(i)}, X, N_s)$  /*
/* Compute the conditional element */
6  $\widehat{\text{T-E}}_A \leftarrow \sum_{l=1}^N \left( \frac{1}{N_s-1} \sum_{i=1}^{N_s} \left[ \mathbb{1}_{\mathcal{F}_t}(\widetilde{X}_l^{A,(i)}) - \frac{1}{N_s} \sum_{h=1}^{N_s} \mathbb{1}_{\mathcal{F}_t}(\widetilde{X}_l^{A,(h)}) \right]^2 \right) \times (N\widehat{p}_t^Y(1 - \widehat{p}_t^Y))^{-1}$  /*
/* Aggregation step */
7 for  $j = 1, \dots, d$  do /*
8  $\widehat{\text{T-Sh}}_{j,\text{MC}} \leftarrow 0$ 
9 for  $A \subset \{-j\}$  do /*
/* Apply the Shapley weights to every computed increments */
10  $\widehat{\text{T-Sh}}_{j,\text{MC}+} = \frac{1}{d} \binom{d-1}{|A|}^{-1} (\widehat{\text{T-E}}_{A \cup \{j\}} - \widehat{\text{T-E}}_A)$ 

```

---

## 5. Analytical results using a linear model with Gaussian inputs

To illustrate the behavior of the proposed indices, a first toy-case involving a linear model and multivariate Gaussian inputs is presented. The main advantage is that analytical results can be derived for the marginal distributions of all subsets of inputs, their conditional distribution, and the distribution of the output given a subset of inputs. Moreover, analytical formulas can be obtained for both the target Sobol' indices and the target Shapley effects.

Let  $(\beta_0, \beta) = (\beta_0, \beta_1, \dots, \beta_d)^\top \in \mathbb{R}^{d+1}$ ,  $\mu = (\mu_1, \dots, \mu_d)^\top \in \mathbb{R}^d$  and  $\Sigma \in \mathcal{M}_d(\mathbb{R})$  a full-rank symmetric  $(d \times d)$  square matrix. Assume that  $X \sim \mathcal{N}_d(\mu, \Sigma)$ , and that the model output writes

$$Y = \beta_0 + \beta^\top X. \quad (34)$$

Then, one has  $Y \sim \mathcal{N}(\beta_0 + \beta^\top \mu, \beta^\top \Sigma \beta)$  and, for any  $A \in \mathcal{P}_d$ ,  $(Y|X_A = x_A) \sim \mathcal{N}(\widetilde{\mu}_A, \widetilde{\Sigma}_A)$  with

$$\widetilde{\mu}_A = \beta_0 + \beta_A^\top x_A + \beta_{\overline{A}}^\top (\mu_{\overline{A}} + \Sigma_{A,12} \Sigma_{A,22}^{-1} (x_A - \mu_A)), \quad \widetilde{\Sigma}_A = \beta_A^\top (\Sigma_{A,11} - \Sigma_{A,12} \Sigma_{A,22}^{-1} \Sigma_{A,21}) \beta_{\overline{A}}.$$

Moreover, one also needs to recall that

$$(X_{\overline{A}}, X_A)^\top \sim \mathcal{N}_d \left( \begin{pmatrix} \mu_{\overline{A}} \\ \mu_A \end{pmatrix}, \Sigma_A = \begin{pmatrix} \Sigma_{A,11} & \Sigma_{A,12} \\ \Sigma_{A,21} & \Sigma_{A,22} \end{pmatrix} \right)$$

with the partitions of  $\Sigma_A$  having sizes  $\begin{pmatrix} (d - |A|) \times (d - |A|) & (d - |A|) \times |A| \\ |A| \times (d - |A|) & |A| \times |A| \end{pmatrix}$ . These results will lead to some numerical approximations of the theoretical values of  $\text{T-Sh}_j$  for all  $j = 1, \dots, d$  using standard multidimensional integration tools. Here, the function `adaptIntegrate()` from the `cubature` package of the R software has been used, with a fixed error tolerance set to  $10^{-8}$ . This allows for studying basic scenarios in order to validate the behavior of the target Shapley effects.

In the following, the independent inputs' case is firstly studied w.r.t. the threshold  $t$ . Then, a case involving linear correlation between inputs (driven by the parameter  $\rho$ ) is studied accordingly. Finally, a last case where an exogenous input is added. Such case appears as soon as an extra input is correlated to another one without being explicitly involved in the model  $G(\cdot)$ .

### 5.1. Independent standard Gaussian inputs

The first toy-case can be specified by:

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim \mathcal{N}_3 \left( \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right), \quad Y = \sum_{i=1}^3 X_i. \quad (35)$$

In this scenario, the three inputs are equally important in terms of defining  $Y$ , but they should also be equally important for the variable of interest  $\mathbf{1}_{\mathcal{F}_t}(X)$ , as assessed by the target Sobol' indices defined in Eq. (17).

From [19, 57], one can easily deduce that the first-order (FO) target Sobol' indices are all equal to each other. Thus, one has:

$$\text{T-S}_{\text{FO}} \stackrel{\text{def}}{=} \text{T-S}_1 = \text{T-S}_2 = \text{T-S}_3 = \frac{\mathbb{V} \left( \Phi \left( \frac{t-X}{\sqrt{2}} \right) \right)}{\mathbb{V}(\mathbf{1}_{\mathcal{F}_t}(X))}, \quad (36)$$

while the second-order (SO) target Sobol' indices are given by:

$$\text{T-S}_{\text{SO}} \stackrel{\text{def}}{=} \text{T-S}_{\{1,2\}} = \text{T-S}_{\{1,3\}} = \text{T-S}_{\{2,3\}} = \frac{\mathbb{V}(\Phi(t - X'))}{\mathbb{V}(\mathbf{1}_{\mathcal{F}_t}(X))} \quad (37)$$

where  $\Phi(\cdot)$  is the standard Gaussian cumulative distribution function (cdf),  $X \sim \mathcal{N}(0, 1)$  and  $X' \sim \mathcal{N}(0, 2)$ . Finally, one can also show that the third-order (TO) target Sobol' indices are equal to:

$$\text{T-S}_{\text{TO}} \stackrel{\text{def}}{=} \text{T-S}_{\{1,2,3\}} = 1. \quad (38)$$

From Eqs. (36), (37), and (38), and from Property 1, one can deduce that:

$$\text{T-Sh}_1 = \text{T-Sh}_2 = \text{T-Sh}_3 = \frac{1}{3}. \quad (39)$$

Additionally, as the inputs are independent, interpreting the closed target Sobol' indices is meaningful,

and they are equal to:

$$\text{T-S}_{\text{FO}} = \text{T-S}_{\text{FO}} \quad (40)$$

$$\text{T-S}_{\text{SO}} = \text{T-S}_{\text{SO}} - 2\text{T-S}_{\text{FO}} \quad (41)$$

$$\text{T-S}_{\text{TO}} = \text{T-S}_{\text{SO}} - 3(\text{T-S}_{\text{FO}}^{\mathbb{1}} + \text{T-S}_{\text{SO}}^{\mathbb{1}}). \quad (42)$$

These results are illustrated in Figure 3. One can remark that, studying the indicator variable  $\mathbb{1}_{\mathcal{F}_t}(X)$  instead of the model output  $Y$  leads to interaction effects between the inputs, while the target Shapley effects remain constant for all threshold values  $t$ . Such a behavior was expected. However, it highlights the fact the target Shapley effects do not report clearly the interaction effects as target Sobol' indices would do (by considering all the various orders). In some sense, it summarizes the information into a single index which focuses only on correlation effects.

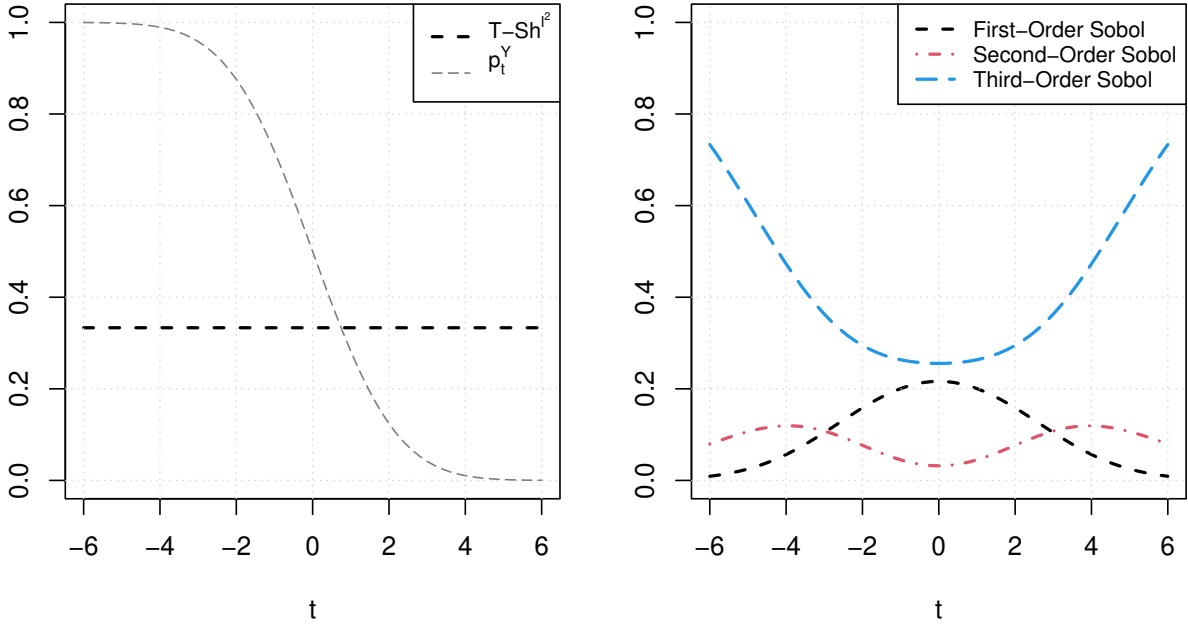


Figure 3: Target Shapley effects (left) and target Sobol' indices (right) for the linear model with standard independent multivariate Gaussian inputs, w.r.t.  $t$ .

### 5.2. Correlated Gaussian inputs with unit variance

The behavior of the target Shapley effects are now studied when a linear dependence is added to the inputs. Since Property 1 still holds without any condition on the dependence structure on the input variables, these indices remain relevant while conveying a practical ease of interpretation. The following model is used:

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} \sim \mathcal{N}_3 \left( \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & \rho \\ 0 & \rho & 1 \end{pmatrix} \right), \quad Y = \sum_{i=1}^3 X_i. \quad (43)$$

where  $-1 < \rho < 1$ . In this scenario one has

$$\text{T-S}_1 = \frac{\mathbb{V}\left(\Phi\left(\frac{t-X}{\sqrt{2(1+\rho)}}\right)\right)}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))}, \quad (44)$$

$$\text{T-S}_2 = \text{T-S}_3 = \frac{\mathbb{V}\left(\Phi\left(\frac{t-X}{\sqrt{2}}\right)\right)}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))}, \quad (45)$$

$$\text{T-S}_{\{1,2\}} = \text{T-S}_{\{1,3\}} = \frac{\mathbb{V}(\Phi(t-X'))}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))}, \quad (46)$$

$$\text{T-S}_{\{2,3\}} = \frac{\mathbb{V}(\Phi(t-X''))}{\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X))} \quad (47)$$

where  $X' \sim \mathcal{N}(0, 1)$ ,  $X'' \sim \mathcal{N}(0, 2)$  and  $X \sim \mathcal{N}(0, 2(1+\rho))$ .

From these results, one can directly remark that  $\text{T-Sh}_2 = \text{T-Sh}_3$ . Note that the values of the target Shapley effects can also be obtained by combinations of target Sobol' indices (see Eq. (25)). These results are illustrated in Figure 4. For fixed threshold values  $t$ , the target Shapley effects of the correlated inputs  $X_2$  and  $X_3$  increases when  $\rho$  increases. This is an expected behavior since, in this case:

$$\mathbb{V}(\mathbb{1}_{\mathcal{F}_t}(X)) = \Phi\left(\frac{t}{\sqrt{3+2\rho}}\right) \left(1 - \Phi\left(\frac{t}{\sqrt{3+2\rho}}\right)\right), \quad (48)$$

and subsequently, for a fixed  $t$ , the variance of the variable of interest will grow with  $\rho$ , as illustrated in Figure 5. This increase in variance due to the correlation between  $X_2$  and  $X_3$  is then allocated through  $\text{T-Sh}_2$  and  $\text{T-Sh}_3$ , which increase with  $\rho$ . On the other hand,  $\text{T-Sh}_1$  decreases accordingly, to accommodate Property 1.

In Figure 4, the behavior of the indices w.r.t.  $\rho$  can also be clearly seen, as  $\text{T-Sh}_1$  is predominantly above  $\text{T-Sh}_2$  and  $\text{T-Sh}_3$  when  $\rho$  is negative, and below when it is positive. This can be explained by the fact that,  $X_2$  and  $X_3$  cancel each other out when their correlation is negative, thus lowering the value of  $\text{T-S}_{\{2,3\}}$  below  $\text{T-S}_{\{1,2\}}$  and  $\text{T-S}_{\{1,3\}}$ , automatically increasing  $\text{T-Sh}_1$  in accordance to Property 1. In the other hand, for positive values of  $\rho$ ,  $\text{T-S}_{\{2,3\}}$  is higher than  $\text{T-S}_{\{1,2\}}$  and  $\text{T-S}_{\{1,3\}}$ , which in turn corresponds to  $\text{T-Sh}_1$  being lower than  $\text{T-Sh}_2 = \text{T-Sh}_3$ .

### 5.3. Quantifying the importance of an exogenous input in the Gaussian setting

In this use case, inspired by [57], the following model is considered:

$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \sim \mathcal{N}_4 \left( \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & \rho \\ 0 & 0 & 1 & 0 \\ 0 & \rho & 0 & 1 \end{pmatrix} \right), \quad Y = X_1 + 6X_2 + 4X_3 \quad (49)$$

where  $X_4$  is an exogenous input, but correlated to  $X_2$ , which should be the most important variable in terms of variance. The threshold is fixed at  $t = 16$ . This scenario enables to check that the target Shapley effects allow for quantifying the importance of  $X_4$  through its correlation with an endogenous input, even though it does not appear directly in the model. Indeed, from the results given in Figure 6,

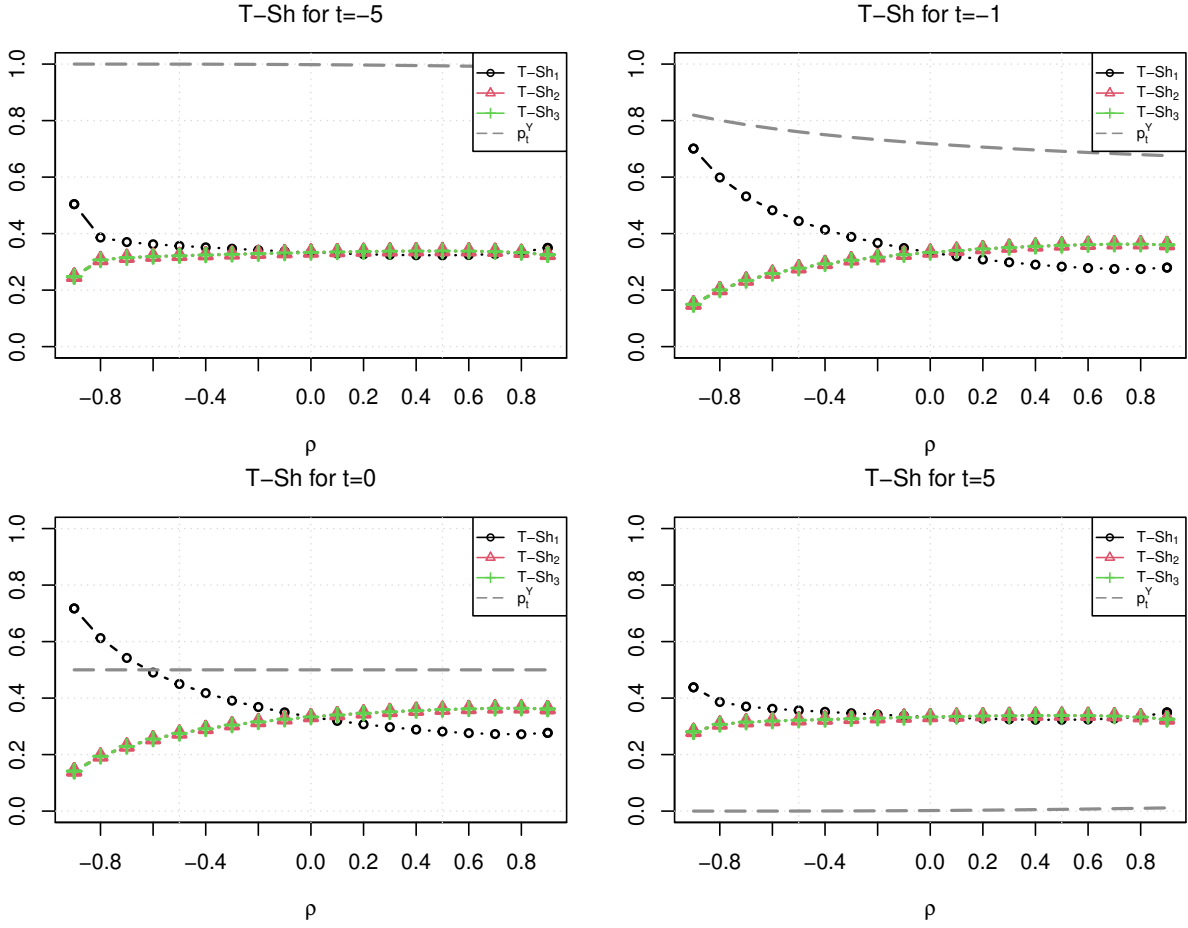


Figure 4: Evolution, w.r.t.  $\rho$  and for various threshold values, of the target Shapley effects of correlated Gaussian standard inputs in a linear model.

one can remark that  $T-Sh_4$  increases when  $\rho$  goes to either 1 or  $-1$ , despite the fact that it is not used directly in the model  $G(\cdot)$ .

## 6. Applications

In this section, two models related to real phenomena and including dependent random inputs are studied in the context of TSA.

### 6.1. A simplified flood model

The target Shapley effects are firstly computed on a simplified model of a river flood [57, 6]. The aim of this model is to study the behavior of the river's water level, by comparison with a fixed dyke height. After a strong simplification of the one-dimensional Saint-Venant equation (with uniform and constant flow rate), the maximal annual water level  $h$  is modeled as

$$h = \left( \frac{Q}{BK_s \sqrt{\frac{Z_m - Z_v}{L}}} \right)^{\frac{3}{2}}, \quad (50)$$

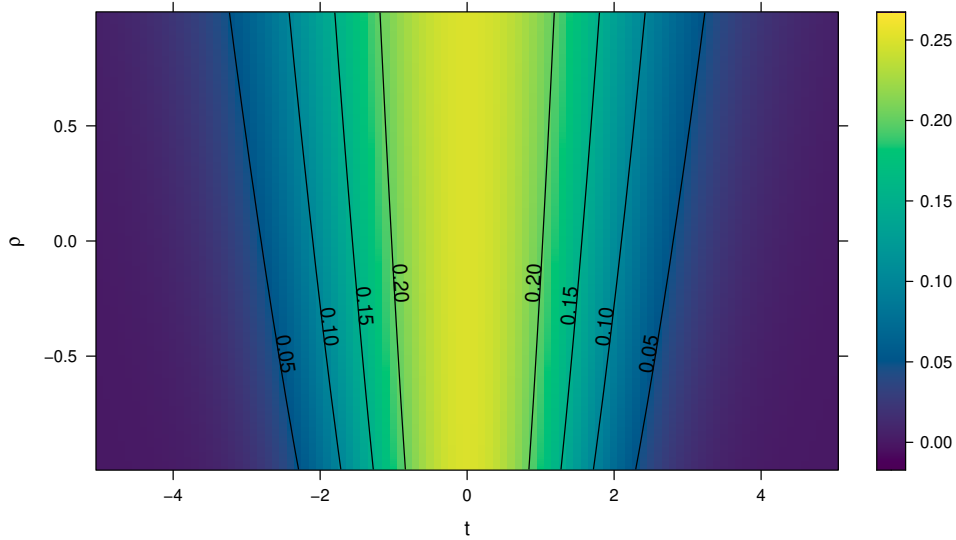


Figure 5: Variance of  $\mathbb{1}_{\mathcal{F}_t}(X)$  w.r.t.  $\rho$  and  $t$  for correlated Gaussian standard inputs with a linear model.

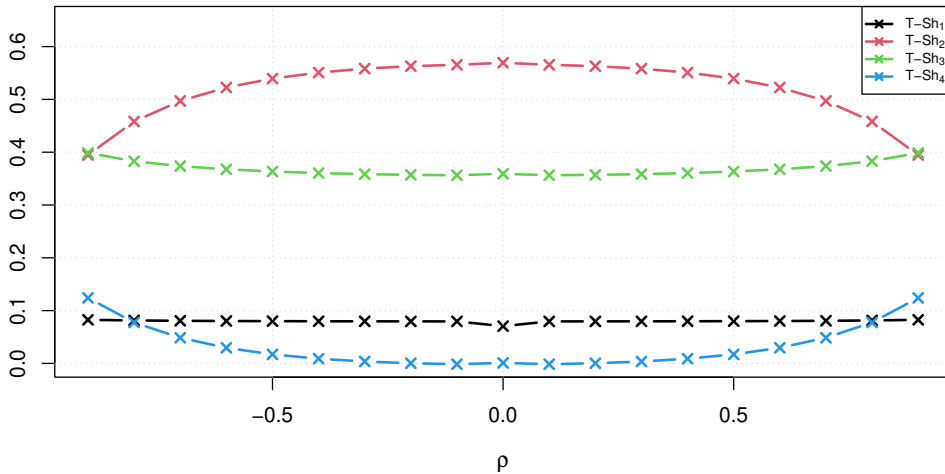


Figure 6: Target Shapley effects for the Gaussian Linear model with an exogenous input, w.r.t.  $\rho$ .

while the model output writes

$$Y = Z_v + h. \quad (51)$$

The inputs' modeling and threshold value are described in Table 1. The problem is of dimension  $d = 6$ . In the present TSA study, the variable of interest is  $\mathbb{1}_{\{G(X) > t\}}(X)$  with  $t$  being the dyke height, fixed to  $t = 54.5$  m. The reference failure probability, computed here with a Monte Carlo sample of large size (here  $10^7$  samples) is equal to  $p_t^Y = 4.5 \times 10^{-3}$ .

As in [24], three pairs of inputs are assumed to be linearly dependent:  $Q$  and  $K_s$  with  $\rho(Q, K_s) = 0.5$ ,  $Z_v$  and  $Z_m$  with  $\rho(Z_v, Z_m) = 0.3$ ,  $L$  and  $B$  with  $\rho(L, B) = 0.3$ . The aim of this use-case is to assess the interest of the target Shapley effects in a complex environment. From [24], it can be shown that, from a GSA standpoint (using a generalized variance decomposition for dependent variables), the two most influential inputs on the annual water level are  $Q$ , the maximal annual flow rate, and

Input	Description	Unit	Distribution
$Q$	maximal annual flow rate	$\text{m}^3.\text{s}^{-1}$	Gumbel(1013, 558) truncated to [500, 3000]
$K_s$	Strickler friction coefficient	-	Normal(30, 7) truncated to [15, $+\infty$ )
$Z_v$	river downstream level	m	Triangular(49, 50, 51)
$Z_m$	river upstream level	m	Triangular(54, 55, 56)
$L$	length of the river stretch	m	Triangular(4990, 5000, 5010)
$B$	river width	m	Triangular(295, 300, 305)
$t$	dyke height (threshold)	m	Fixed to 54.5

Table 1: Input variables and distributions for the flood model.

$Z_v$ , the river downstream level.

A Monte Carlo sample of  $N = 2 \times 10^5$  input realizations is drawn (note that the correlations are injected following the algorithm proposed by [58]). Thus, this procedure leads to the computation of  $N$  model output values. Figure 7 represents the estimated target Shapley effects on the flood case, using the nearest-neighbor procedure depicted in Subsection 4.2 with an arbitrary number of neighbors set at  $N_s = 2$ . Moreover, 300 repetitions of the simulation and the estimation procedure give access to confidence intervals of the estimates (represented by boxplots in Figure 7). From these TSA results, one can notice that  $Q$  is granted an influence of 24.3% ( $\pm 1.3\%$ ),  $K_s$  has 22.6% ( $\pm 1.3\%$ ) and  $Z_v$  around 16.7% ( $\pm 1\%$ ). The other inputs present a share of around 12%. Compared to results obtained by GSA without correlations [6] and with correlations [24], these TSA results with correlations present a much larger effect for  $K_s$  and non-negligible effects for  $Z_m$ ,  $L$  and  $B$ . This was expected due to the interactions induced by the considered TSA variable of interest. This example illustrates the ability of the target Shapley effects to quantify the relative importance of input variables in a complex model presenting statistical dependence between them.

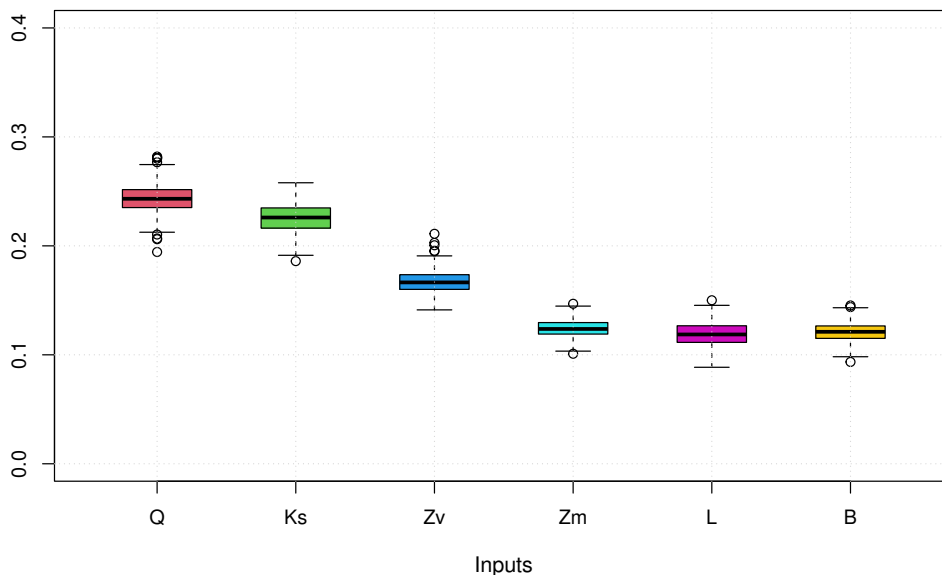


Figure 7: Estimated target Shapley effects for the flood case.

## 6.2. A COVID-19 epidemiological model

In 2020, the COVID-19 crisis has raised major issues in the usefulness of epidemic modeling in order to give useful insights to public policy decision makers. [59] have taken this example to insist on the essential use of GSA on such models, which claim to predict the potential consequences of intervention policies. A first study has been proposed by [60], in the context of COVID-19 in Italy, to assess the sensitivity of model outputs such as quarantined, recovered or dead people to inputs driving intervention policies. Another GSA has been performed in [61] in the French context of the first COVID-19 wave. By using data coming from this last analysis (thanks to the authors' agreement), the goal of this section is to demonstrate how TSA can help to characterize the influence of various uncertain parameters on a real-scale model.

### 6.2.1. Description of the model and its inputs

The deterministic compartmental model developed in [61] (also presented in [62]) is representative of the COVID-19 French epidemic (from March to May) by taking into account the asymptomatic individuals, the testing strategies, the hospitalized individuals, and people going to Intensive Care Unit (ICU). Using several assumptions, it is based on a system of 10 ordinary differential equations which are not developed here for a sake of conciseness (see [62, 61] for more information).

Table 2 presents the 20 input parameters with their prior distribution (chosen from literature studies), which form the inputs  $X$ , assumed to be independent between each other. For the present study, our variable of interest, which is a particular model output, then writes

$$U_{\max}^p = \max_{v \in \text{time range}} \{U_v(X)\} \quad (52)$$

where  $U_v$  is the the number of hospitalized patients in ICU at time  $v$ . Note that the ‘‘p’’ in  $U_{\max}^p$  stands for ‘‘prior’’ as this quantity corresponds to the variable of interest before any calibration w.r.t. the available data.

In [61], after a first screening step allowing to suppress non-influential inputs, the model is calibrated on real data by using a Bayesian calibration technique. After the analysis of this step, the selected remaining inputs are

$$X_{\text{sel}} = (p_a, N_a, N_s, R_0, t_0, \mu, N, I_0^-)^\top \quad (53)$$

and their distributions are obtained from a sample given by the calibration process. The non-influential inputs are fixed to their nominal values and the posterior variable of interest becomes

$$U_{\max} = \max_{v \in \text{time range}} \{U_v(X_{\text{sel}})\} \quad (54)$$

with  $U_{\max}$  being the maximum number of hospitalized people in ICU who need special medical care on the studied temporal range, and  $U_v$  is the number of hospitalized patients in ICU at time  $v$ .

### 6.2.2. Input importance for ICU bed shortage

The central question of this study would be to determine which inputs influence the event of a country experiencing a shortage of ICU bed during the time period. For that purpose, one can introduce a threshold  $t$ , which represents the total number of ICU beds in the country, which is



Input	Description	Prior distribution
$p_a$	Conditioned on being infected, the probability of having light symptoms or no symptoms	$\mathcal{U}(0.5, 0.9)$
$p_H$	Conditioned on being mild/severely ill, the probability of needing hospitalization ( $H$ or $U$ )	$\mathcal{U}(0.15, 0.2)$
$p_U$	Conditioned on going to hospital, the probability of needing ICU	$\mathcal{U}(0.15, 0.2)$
$p_{HD}$	Conditioned on being hospitalized but not in ICU, the probability of dying	$\mathcal{U}(0.15, 0.25)$
$p_{UD}$	Conditioned on being admitted to ICU, the probability of dying	$\mathcal{U}(0.2, 0.3)$
$N_a$	If asymptomatic, number of days until recovery	$\mathcal{U}(8, 12)$
$N_s$	If symptomatic, number of days until recovery without hospital	$\mathcal{U}(8, 12)$
$N_{IH}$	If severe symptomatic, number of days until hospitalization	$\mathcal{U}(8, 12)$
$N_{HD}$	If in $H$ , number of days until death	$\mathcal{U}(15, 20)$
$N_{UD}$	If in ICU, number of days until death	$\mathcal{U}(8, 12)$
$N_{HR}$	If hospitalized but not in ICU, the number of days until recovery	$\mathcal{U}(15, 25)$
$N_{UR}$	If in ICU, number of days until recovery	$\mathcal{U}(15, 25)$
$R_0$	Basic reproducing number	$\mathcal{U}(3, 3.5)$
$t_0$	Starting date of epidemics (in 2020)	$\mathcal{U}(01/25, 02/24)$
$\mu$	Decaying rate for transmission (after social distancing and lockdown)	$\mathcal{U}(0.03, 0.08)$
$N$	Date of effect of social distancing and lockdown	$\mathcal{U}(20, 50)$
$\lambda_1$	Type-1 testing rate	$\mathcal{U}(1e-4, 1e-3)$
$p_{HU}$	Conditioned on being hospitalized in $H$ , the probability of needing ICU	$\mathcal{U}(0.15, 0.2)$
$N_{HU}$	If in $H$ , number of days until ICU	$\mathcal{U}(1, 10)$
$I_0^-$	Number of infected undetected at the start of epidemics	$\mathcal{U}(1, 100)$

Table 2: Model inputs and their prior distribution.  $H$  is the number of hospitalized individuals with severe symptoms.  $U$  is the number of hospitalized individuals in ICU.

assumed to be constant during the studied time period. The new variable of interest would then be  $\mathbb{1}_{\{U_{\max}^P > t\}}(X)$  for the full compartmental model (preliminary study) and  $\mathbb{1}_{\{U_{\max} > t\}}(X_{\text{sel}})$  for the model with selected inputs (post-calibration study). Two input-output samples of size  $n = 5000$  are available. The first one (preliminary study) includes all the inputs following their prior distribution (see Table 2) and the corresponding output  $U_{\max}^P$  of the compartmental model. The second one (post-calibration study) is composed of a sample of  $X_{\text{sel}}$  after the Bayesian calibration, and the corresponding output  $U_{\max}$  of the compartmental model with the non-selected inputs fixed to their nominal values.

Five different thresholds are studied on  $U_{\max}^P$ :  $5 \cdot 10^3$ ,  $10^4$ ,  $5 \cdot 10^4$ ,  $10^5$  and  $2 \cdot 10^5$ , with respectively 58.1%, 7.7%, 22%, 10.1% and 2.2% of the total output samples being in a failure state. This allows to illustrate the behavior of the target Shapley effects when the failure probability decreases. The threshold of 6300 has been chosen for  $U_{\max}$ , with 10.9% of the total output samples being above this threshold. Figure 8 illustrates two different thresholds, and the corresponding estimated failure probability on the histogram of both outputs.

The target Shapley effects have been estimated using a variant of the estimation scheme presented in Subsection 4.2, with a fixed number of random permutations of  $10^3$ , and with a number of neighbors set to 3, following the rule of thumb guideline of [36], due to the sheer complexity of this estimation algorithm. Since the compartmental model is deterministic, the target Shapley effects have been forced

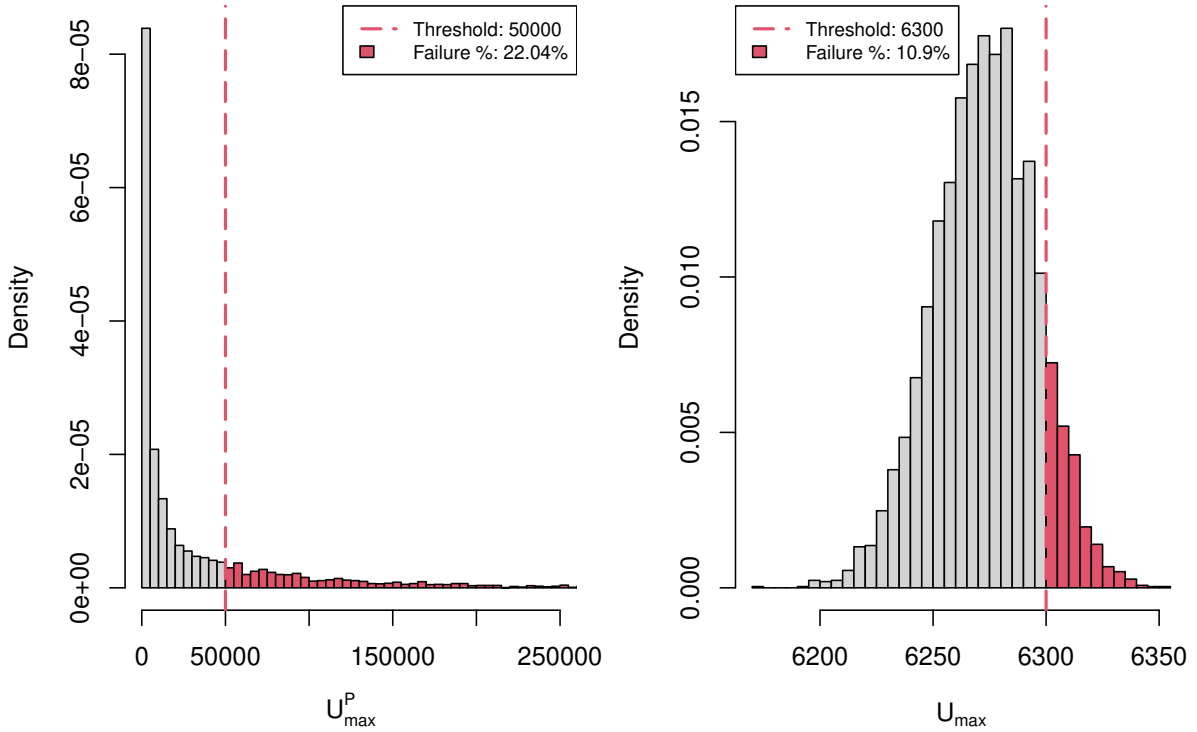


Figure 8: Illustration of thresholds on the histograms of  $U_{\max}^P$  (left) and  $U_{\max}$  (right).

to sum up to one. Figure 9 presents the main results for  $U_{\max}^P$ , with the red dotted line being the average influence of an input, in the case of similar importance (i.e.,  $\frac{1}{20}$ ). One can remark that for less restrictive thresholds (i.e., threshold for which the failure probability is high), the input  $N$ , the effective date of lockdown/social distancing measures, seem to be the most influential, reaching more than 50% of the TSA variable of interest's variance. However, as soon as the threshold becomes more and more restrictive (i.e., the failure probability becomes lower and lower), the effect of  $N$  decreases, and the effects of the other inputs increase accordingly, in order to reach what seem to be an equilibrium at the value  $\frac{1}{20}$ . This behavior can be explained by two main reasons:

- As outlined in Subsection 3.1, the nature of a restrictive TSA variable of interest induces high interaction between the inputs;
- The Shapley allocation system, when applied to variance as a production value, redistributes the interaction effects equally between all inputs (there is no correlation between inputs in this prior study).

One can argue that, as soon as  $t$  becomes very restrictive, the combined interaction effects outweighs the effect of  $N$  itself, and since these effects are equitably distributed among all the inputs, their share will tend to go towards  $\frac{1}{20}$ .

For the post-calibration study, some selected inputs  $X_{\text{sel}}$  are linearly correlated (see Figure 10 - top left). This is typically the case for  $N$  and  $\mu$ , with an estimated correlation coefficient  $\hat{\rho}(N, \mu) = 0.69$ , and for  $R_0$  and  $N$  with an estimated correlation coefficient of  $\hat{\rho}(N, R_0) = -0.66$ . This correlation

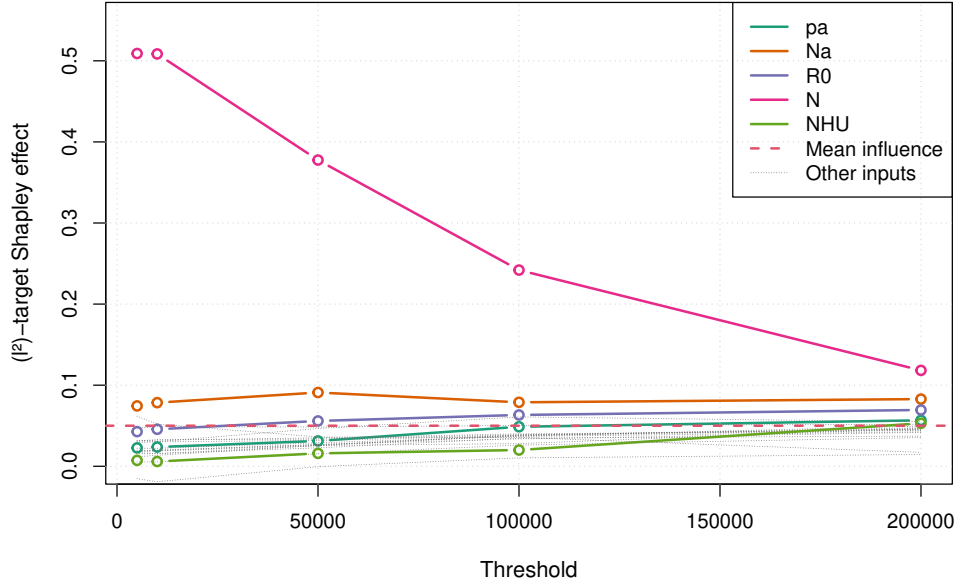


Figure 9: Target Shapley effects for  $\mathbb{1}_{\{U_{\max}^P > \epsilon\}}(X)$  for different thresholds.

structure does not allow to get interpretable Sobol' indices, as outlined in Section 2, which motivates the use of Shapley-inspired indices. The Shapley effects and the target Shapley effects of  $X_{\text{sel}}$  for  $U_{\max}$  have been computed using the nearest-neighbor procedure, with a fixed number of neighbors of 3, and forced to sum to one because of the deterministic nature of the model.

One can remark that  $N_a$ , the number of days until recovery, seem to be the most important input in explaining the number maximum number of ICU patients on the studied time range, with a Shapley effect of around 35% of the output variance. The inputs  $p_a$ ,  $N_s$ ,  $R_0$  and  $N$  seem to present average effects, that is around  $\frac{1}{8}$ , while  $t_0$ ,  $\mu$  and  $I_0^-$  seem to be less influential, with around 5% of explained variance each.

However, focusing on the event of a ICU bed shortage, one can remark that the target Shapley effect of  $N_a$  is lower (around 22%), with the influence of  $N$  being higher (around 15%) than their Shapley effects. Moreover,  $t_0$ ,  $\mu$  and  $I_0^-$  present higher TSA effects, i.e., slightly under 10%, due to the interaction induced by the indicator function. One can also remark that the influence of  $N_s$  is higher than the one of  $R_0$  in the TSA setting, which was the other way around for the Shapley effects. This would indicate that  $N_s$ , the number of days until recovery for a symptomatic patient without hospitalization, influences more the event of a bed shortage than the basic reproducing number of the virus,  $R_0$ .

## 7. Conclusion

This paper aims at proposing a set of new indices adapted to target sensitivity analysis while being able to handle correlated inputs. The goal is to quantify the importance of inputs on the occurrence of critical failure event of the system under study. The proposed indices are based on a cooperative Shapley procedure which aims at allocating the effects of the interaction and correlation

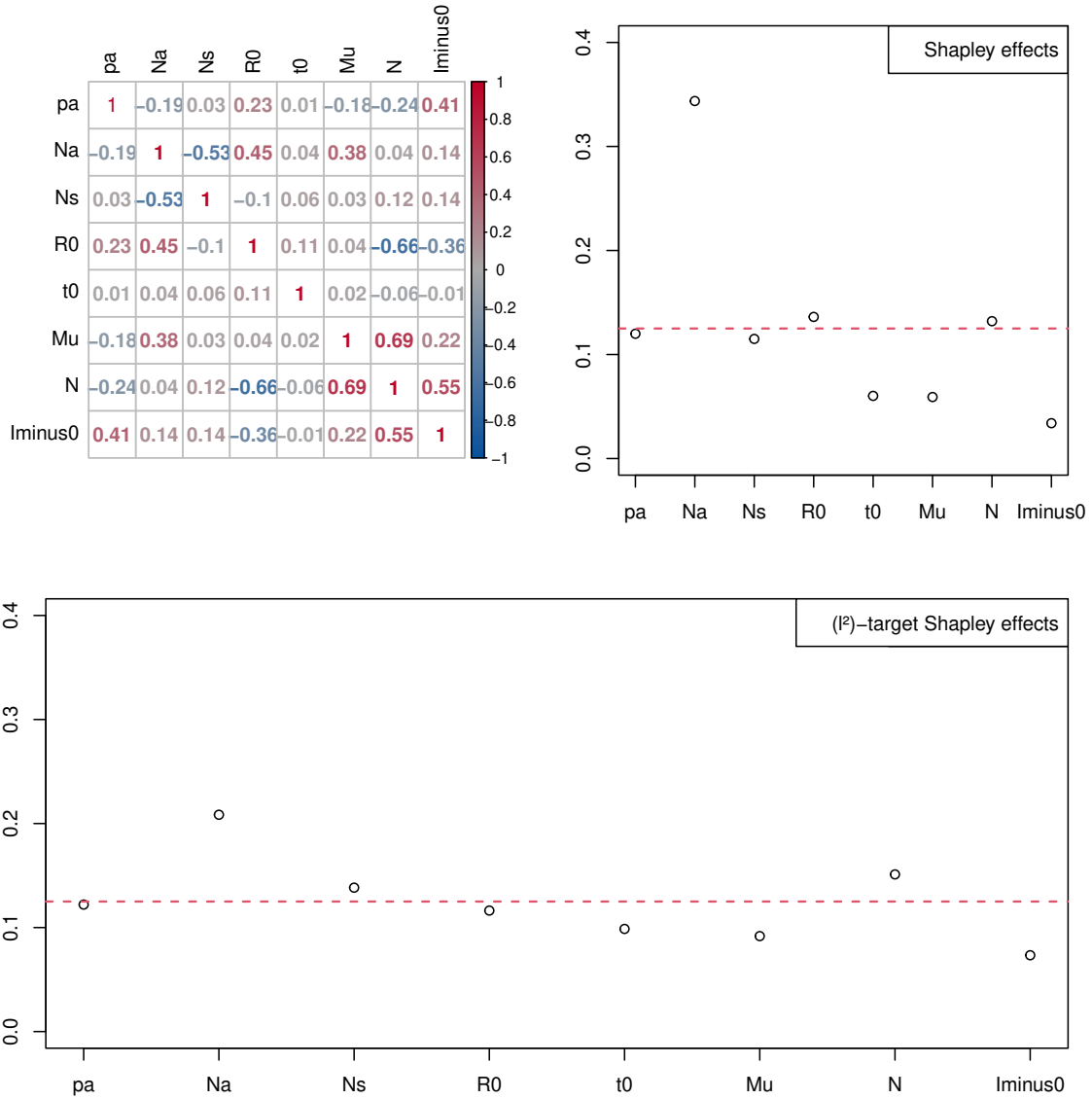


Figure 10: Input correlation matrix (top left), Shapley effects for  $U_{\max}$  (top right) and target Shapley effects (bottom) for  $\mathbb{1}_{\{U_{\max} > t\}}(X_{\text{sel}})$ .

equally between all the inputs in the same manner as done by the Shapley effects in global sensitivity analysis. Thus, a general class of distance-based indices is proposed, namely the  $(\mathcal{D})$ -target Shapley effects and some relevant properties are highlighted. Depending on the choice of the distance  $\mathcal{D}$ , well-known preexisting indices can be used as cost functions in the Shapley formulation. Therefore, these indices allow for the allocation, among the different inputs, of shares of several dispersion statistics (e.g., mean absolute deviation for the  $\ell^1$  case, variance for the  $\ell^2$  case). This produces easily usable indices in practice, as they can directly be interpreted as percentage of the dispersion statistic, allocated to each input. This versatile procedure allows to produce input importance measures according to a

specific metric, driven by the choice of the distance.

In particular, the ( $\ell^2$ )-target Shapley effects (called target Shapley effects to simplify), which represents percentages of variance, have been studied more intensively and two dedicated estimation methods have been proposed. These particular indices have then been applied, analyzed and discussed through simple Gaussian toy-cases. Finally, two real-world use-cases have been studied: the modeling of a river flood and the problem of ICU bed shortage during the COVID-19 pandemic. These indices reveal to be able to detect influential inputs in the context of correlated inputs. For target sensitivity analysis, such a tool is valuable and can be used as a complement of more standard procedures. The clear advantage is that only one set of indices is required by an analyst in order to produce easily interpretable and meaningful insights. Moreover, the proposed indices can be estimated in the context of given-data which can be adapted to real applications for which no computer model is available. However, the major drawbacks of the approach are mainly related to the target aspect of the analysis. Indeed, as soon as the event gets more and more rare, all the inputs tend to be influential and making a clear distinction between interactions and correlation effects becomes difficult.

A first perspective could be to improve the estimation strategies. The sampling-based method could benefit from a better sampling scheme, such as importance sampling, as described in [63], which could reduce the estimator's variance. Recent results from [64] using copulas also shows promising tracks to have efficient estimation of Shapley effects. Moreover, adapting recent results from [35], with a link between the target Sobol' indices and the Squared Mutual Information [65], should allow for other possibilities of given-data estimation methods. Another method based on a random forest given-data procedure, explored in [66] in the context of quantile-oriented importance measure estimation, could also yield promising results if transposed to a reliability-oriented setting.

Finally, even if the Shapley attribution system allows to deal with input statistical dependencies, a finer decomposition is lacking in order to quantify the origin of each effect (e.g., statistical dependence and interaction). Further works should use the recent developments in [67] in order to quantify interaction effects, with a transposition to the target sensitivity analysis setting. Finally, it has been shown in [68] that the Shapley attribution system is equivalent to a maximum entropy distribution (e.g., uniform) over all possible orderings of inputs (the Shapley weights). Developments towards other forms of allocation systems where the allocation is driven from data could also open a path for further developments.

## Acknowledgments

We are grateful to as Sébastien Da Veiga, Clémentine Prieur and Fabrice Gamboa for helpful discussions and for having provided the dataset on the COVID-19 model.

## References

- [1] E. De Rocquigny, N. Devictor, S. Tarantola, *Uncertainty in industrial practice: a guide to quantitative uncertainty management*, Wiley, 2008.
- [2] K. Beven, *Environmental Modelling: An Uncertain Future?*, CRC Press, 2008.

- [3] F. Pianosi, K. Beven, J. Freer, J. Hall, J. Rougier, D. Stephenson, T. Wagener, Sensitivity analysis of environmental models: A systematic review with practical workflow, *Environmental Modelling & Software* 79 (2016) 214–232.
- [4] S. Razavi, A. Jakeman, A. Saltelli, C. Prieur, B. Iooss, E. Borgonovo, E. Plischke, S. Lo Piano, T. Iwanaga, W. Becker, S. Tarantola, J. Guillaume, J. Jakeman, H. Gupta, N. Melillo, G. Rabiti, V. Chabridon, Q. Duan, X. Sun, S. Smith, R. Sheikholeslami, N. Hosseini, M. Asadzadeh, A. Puy, S. Kucherenko, H. Maier, The future of sensitivity analysis: An essential discipline for systems modelling and policy making, *Environmental Modelling and Software*, In press (104954) (2020).
- [5] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, S. Tarantola, *Global Sensitivity Analysis. The Primer*, Wiley, 2008.
- [6] B. Iooss, P. Lemaître, A review on global sensitivity analysis methods, in: C. Meloni, G. Dellino (Eds.), *Uncertainty management in Simulation-Optimization of Complex Systems: Algorithms and Applications*, Springer, 2015, pp. 101–122.
- [7] M. Lemaire, A. Chateauneuf, J.-C. Mitteau, *Structural Reliability*, ISTE Ltd & John Wiley & Sons, Inc., 2009.
- [8] Y. Richet, V. Bacchi, Inversion algorithm for civil flood defense optimization: Application to two-dimensional numerical model of the garonne river in france, *Frontiers in Environmental Science* 7 (2019) 160.
- [9] R. T. Rockafellar, J. O. Royset, Engineering Decisions under Risk Averseness, *ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering* 1 (2) (2015) 1–12.
- [10] J. Morio, M. Balesdent, *Estimation of Rare Event Probabilities in Complex Aerospace and Other Systems: A Practical Approach*, Woodhead Publishing, Elsevier, 2015.
- [11] Y.-T. Wu, Computational Methods for Efficient Structural Reliability and Reliability Sensitivity Analysis, *AIAA Journal* 32 (8) (1994) 1717–1723.
- [12] S. Song, Z. Lu, H. Qiao, Subset simulation for structural reliability sensitivity analysis, *Reliability Engineering and System Safety* 94 (2) (2009) 658–665.
- [13] P. Wei, Z. Lu, W. Hao, J. Feng, B. Wang, Efficient sampling methods for global reliability sensitivity analysis, *Computer Physics Communications* 183 (2012) 1728–1743.
- [14] V. Chabridon, Reliability-oriented sensitivity analysis under probabilistic model uncertainty – Application to aerospace systems, Ph.D. thesis, Université Clermont Auvergne (2018).
- [15] G. Perrin, G. Defaux, Efficient Evaluation of Reliability-Oriented Sensitivity Indices, *Journal of Scientific Computing* (2019) 1–23.
- [16] I. M. Sobol, Sensitivity estimates for nonlinear mathematical models, *Mathematical Modelling and Computational Experiments* 1 (1993) 407–414.

- [17] E. Borgonovo, A new uncertainty importance measure, *Reliability Engineering & System Safety* 92 (6) (2007) 771–784.
- [18] H. Raguét, A. Marrel, Target and conditional sensitivity analysis with emphasis on dependence measures, Working paper (2018, URL <https://arxiv.org/abs/1801.10047>).
- [19] L. Li, Z. Lu, F. Jun, W. Bintuan, Moment-independent importance measure of basic variable and its state dependent parameter solution, *Structural Safety* 38 (2012) 40–47.
- [20] A. Marrel, V. Chabridon, Statistical developments for target and conditional sensitivity analysis: application on safety studies for nuclear reactor, Preprint HAL, hal-02541142v2 (2020, URL <https://hal.archives-ouvertes.fr/hal-02541142v2>).
- [21] W. Hoeffding, A class of statistics with asymptotically normal distribution, *The Annals of Mathematical Statistics* 19 (3) (1948) 293–325.
- [22] J. Jacques, C. Lavergne, N. Devictor, Sensitivity analysis in presence of model uncertainty and correlated inputs, *Reliability Engineering & System Safety* 91 (2006) 1126–1134.
- [23] G. Li, H. Rabitz, P. Yelvington, O. Oluwole, F. Bacon, C. Kolb, J. Schoendorf, Global sensitivity analysis for systems with independent and/or correlated inputs, *Journal of Physical Chemistry* 114 (2010) 6022–6032.
- [24] G. Chastaing, F. Gamboa, C. Prieur, Generalized Hoeffding-Sobol decomposition for dependent variables - Application to sensitivity analysis, *Electronic Journal of Statistics* 6 (2012) 2420–2448.
- [25] C. Xu, G. Z. Gertner, Uncertainty and sensitivity analysis for models with correlated parameters, *Reliability Engineering & System Safety* 93 (2008) 1563–1573.
- [26] T. Mara, S. Tarantola, Variance-based sensitivity indices for models with dependent inputs, *Reliability Engineering & System Safety* 107 (2012) 115–121.
- [27] T. Mara, S. Tarantola, P. Annoni, Non-parametric methods for global sensitivity analysis of model output with dependent inputs, *Environmental Modeling & Software* 72 (2015) 173–183.
- [28] N. Benoumechiara, K. Elie-Dit-Cosaque, Shapley effects for sensitivity analysis with dependent inputs: bootstrap and kriging-based algorithms, *ESAIM: Proceedings and Surveys* 65 (2019) 266–293.
- [29] N. Do, S. Razavi, Correlation effects? A major but often neglected component in sensitivity and uncertainty analysis, *Water Resources Research* 56 (e2019WR025436) (2020).
- [30] L. S. Shapley, A value for n-person games, in: H. Kuhn, A. W. Tucker (Eds.), *Contributions to the Theory of Games, Volume II*, *Annals of Mathematics Studies*, Princeton University Press, Princeton, NJ, 1953, Ch. 17, pp. 307–317.
- [31] M. Osborne, A. Rubinstein, *A Course in Game Theory*, MIT Press, 1994.

- [32] A. B. Owen, Sobol' indices and Shapley value, *SIAM/ASA Journal of Uncertainty Quantification* 2 (2014) 245–251.
- [33] A. B. Owen, C. Prieur, On Shapley value for measuring importance of dependent inputs, *SIAM/ASA Journal of Uncertainty Quantification* 5 (2017) 986–1002.
- [34] B. Iooss, C. Prieur, Shapley effects for Sensitivity Analysis with correlated inputs : Comparisons with Sobol' Indices, *Numerical Estimation and Applications, International Journal for Uncertainty Quantification* 9 (5) (2019) 493–514.
- [35] A. Spagnol, Kernel-based sensitivity indices for high-dimensional optimization problems, Ph.D. thesis, Ecole des Mines de Saint-Etienne (2020).
- [36] B. Broto, F. Bachoc, M. Depecker, Variance reduction for estimation of shapley effects and adaptation to unknown input distribution, *SIAM/ASA Journal on Uncertainty Quantification* 8 (2) (2020) 693–716.
- [37] R. Christensen, *Linear models for multivariate, time series and spatial data*, Springer-Verlag, 1990.
- [38] T. Hastie, R. Tibshirani, J. Friedman, *The elements of statistical learning*, Springer, 2002.
- [39] J. Helton, J. Johnson, C. Salaberry, C. Storlie, Survey of sampling-based methods for uncertainty and sensitivity analysis, *Reliability Engineering & System Safety* 91 (2006) 1175–1209.
- [40] J. Johnson, J. LeBreton, History and use of relative importance indices in organizational research, *Organizational Research Methods* 7 (2004) 238–257.
- [41] L. Clouvel, Uncertainty quantification of the fast flux calculation for a PWR vessel, Ph.D. thesis, Université Paris-Saclay (2019).
- [42] R. H. Lindeman, P. F. Merenda, R. Z. Gold, *Introduction to bivariate and multivariate analysis*, Scott Foresman and Company, Glenview, IL, 1980.
- [43] U. Grömping, Relative importance for linear regression in R: the Package `relaimpo`, *Journal of Statistical Software* 17 (2006) 1–27.
- [44] J. Nossent, P. Elsen, W. Bauwens, Sobol' sensitivity analysis of a complex environmental model, *Environmental Modelling & Software* 26 (12) (2011) 1515 – 1525.
- [45] E. Song, B. Nelson, J. Staum, Shapley effects for global sensitivity analysis: Theory and computation, *SIAM/ASA Journal on Uncertainty Quantification* 4 (1) (2016) 1060–1083.
- [46] P. Derennes, J. Morio, F. Simatos, Simultaneous estimation of complementary moment independent and reliability-oriented sensitivity measures, *Mathematics and Computers in Simulation* 182 (2021) 721–737.
- [47] J. Morio, Extreme quantile estimation with nonparametric adaptive importance sampling, *Simulation Modelling Practice and Theory* 27 (2012) 76–89.



- [48] V. Chabridon, M. Balesdent, G. Perrin, J. Morio, J.-M. Bourinet, N. Gayton, *Mechanical Engineering Under Uncertainties*, Wiley - ISTE Ltd, 2020, Ch. ‘Global reliability-oriented sensitivity analysis under distribution parameter uncertainty’, pp. 1–43.
- [49] L. Cui, Z. Lu, X. Zhao, Moment-independent importance measure of basic random variable and its probability density evolution solution, *Science China Technical Sciences* 53 (10) (2010) 1138–1145.
- [50] J.-C. Fort, T. Klein, N. Rachdi, New sensitivity analysis subordinated to a contrast, *Communications in Statistics - Theory and Methods* 45 (15) (2016) 4349–4364.
- [51] T. Browne, J.-C. Fort, B. Iooss, L. Le Gratiet, Estimate of quantile-oriented sensitivity indices, HAL, hal-01450891, version 1 (2017).
- [52] V. Maume-Deschamps, I. Niang, Estimation of quantile oriented sensitivity indices, *Statistics and Probability Letters* 134 (2018) 122–127.
- [53] S. Kucherenko, S. Song, L. Wang, Quantile based global sensitivity measures, *Reliability Engineering and System Safety* 185 (2019) 35–48.
- [54] L. Li, I. Papaioannou, D. Straub, Global reliability sensitivity estimation based on failure samples, *Structural Safety* 81 (2019) 101871.
- [55] L. Li, Z. Lu, C. Chen, Moment-independent importance measure of correlated input variable and its state dependent parameter solution, *Aerospace Science and Technology* 48 (2016) 281–290.
- [56] B. Iooss, S. Da Veiga, A. Janon, G. Pujol, sensitivity: Global Sensitivity Analysis of Model Outputs, R package version 1.23.1 (2020).  
URL <https://CRAN.R-project.org/package=sensitivity>
- [57] P. Lemaitre, *Analyse de sensibilité en fiabilité des structures* (in English), Ph.D. thesis, Université de Bordeaux (2014).
- [58] E. Schumann, Generating correlated uniform variates. (2009).  
URL <http://comisef.wikidot.com/tutorial:correlateduniformvariates>
- [59] A. Saltelli, G. Bammer, I. Bruno, E. Charters, M. D. Fiore, et al., Five ways to ensure that models serve society: a manifesto (short comments), *Nature* 582 (2020) 482–484.
- [60] X. Lu, E. Borgonovo, Is time to intervention in the COVID-19 outbreak really important? A global sensitivity analysis approach, Preprint (2020, arXiv:2005.01833).
- [61] S. Da Veiga, F. Gamboa, C. Prieur, B. Iooss, Basics and trends in sensitivity analysis, SIAM, In revision, 2021.
- [62] S. Da Veiga, Calibration and sensitivity analysis of a COVID-19 epidemics model, Meeting AppliBUGS (Applications du Bayesian Unified Group of Statisticians), December 2020.  
URL <http://genome.jouy.inra.fr/applibugs/applibugs.rencontres.html>

- [63] R. Y. Rubinstein, D. P. Kroese, Simulation and the Monte Carlo method, Second Edition, Wiley, 2008.
- [64] G. Sarazin, P. Derennes, J. Morio, Estimation of high-order moment-independent importance measures for Shapley value analysis, Applied Mathematical Modelling 88 (2020) 396–417.
- [65] M. Sugiyama, Machine learning with squared-loss mutual information, Entropy 15 (1) (2012) 80–112.
- [66] K. Elie-Dit-Cosaque, Développement de mesures d’incertitudes pour le risque de modèle dans des contextes incluant de la dépendance stochastique (in English), Ph.D. thesis, Université Claude Bernard Lyon 1 (2020).
- [67] G. Rabitti, E. Borgonovo, A Shapley–Owen index for interaction quantification, SIAM/ASA Journal on Uncertainty Quantification 7 (3) (2019) 1060–1075.
- [68] E. S. Soofi, J. J. Retzer, M. Yasai-Ardekani, A framework for measuring the importance of variables with applications to management research and decision models, Decision Sciences 31 (2000) 595 – 625.

## Appendix A. ANOVA and Sobol’ indices

In the general non-linear case, as for the ANOVA of the linear model case (see Subsection 2.1), the idea is to find a general decomposition of the output variance. This can be done through the decomposition of a function with finite variance ( $L^2$  mathematical property), called the *Hoeffding decomposition* [21], which allows to rewrite  $G(X)$  as a sum of centered components related to each possible subset of inputs. For example, in the case of a model with three inputs  $X = (X_1, X_2, X_3)$ ,  $G(X)$  can be decomposed into four components:

$$\begin{aligned}
 G(X) &= G_\emptyset && \text{(Mean behavior)} \\
 &+ G_1(X_1) + G_2(X_2) + G_3(X_3) && \text{(First-order)} \\
 &+ G_{\{1,2\}}(X_1, X_2) + G_{\{1,3\}}(X_1, X_3) + G_{\{2,3\}}(X_2, X_3) && \text{(Second-order)} \\
 &+ G_{\{1,2,3\}}(X). && \text{(Third-order)}
 \end{aligned}$$

Moreover, if the inputs are assumed to be independent, each term is orthogonal to one another and writes

$$G_A(x_A) = \sum_{B \subset A} (-1)^{|A|-|B|} \mathbb{E}[G(X)|X_B = x_B] \tag{A.1}$$

where  $A \in \mathcal{P}_d$  is a subset of indices and  $\mathcal{P}_d$  the set of all possible subsets of  $\{1, \dots, d\}$ ,  $|A|$  is the cardinal of  $A$  and  $X_A$  denotes the subset of inputs, selected by the indices in  $A$  ( $X_A = (X_i)_{i \in A}$ ). Then, the Hoeffding decomposition is unique and leads to a variance decomposition called “functional ANOVA”:

$$\mathbb{V}[G(X)] = \sum_{A \in \mathcal{P}_d, A \neq \emptyset} \mathbb{V}[G_A(x_A)]. \tag{A.2}$$

This leads to the definition of the Sobol' indices [16]:

$$S_A = \frac{\mathbb{V}[G_A(X_A)]}{\mathbb{V}[G(X)]} = \frac{\sum_{B \subset A} (-1)^{|A|-|B|} \mathbb{V}(\mathbb{E}[G(X) | X_B])}{\mathbb{V}[G(X)]}. \quad (\text{A.3})$$

The sum of the Sobol' indices over all subset on inputs  $A \in \mathcal{P}_d$  being equal to one, they can be directly interpreted as the percentage of the output variance due to each subset of input [16, 5]. The Sobol' indices of higher orders than one can be interpreted as a means of quantifying the share of variance due to the interaction effects induced by the structure of the model  $G(\cdot)$  between the selected subset of inputs.

Another useful sensitivity index is the closed Sobol' index [16] which writes

$$S_A^{\text{clos}} = \sum_{B \subset A} S_B = \frac{\mathbb{V}(\mathbb{E}[G(X) | X_A])}{\mathbb{V}[G(X)]} \quad (\text{A.4})$$

In the independent setting, it can be interpreted as the percentage of variability induced by all the variables in a selected subset and their interactions. Figure A.11 provides an illustration of the Sobol' indices and the closed Sobol' indices for a model with three inputs. Each Venn diagram represents the variance of the output, with the representation of each of the two Sobol' indices presented above. While this representation is useful in the GSA context, it relies on the assumption of independence between the inputs.

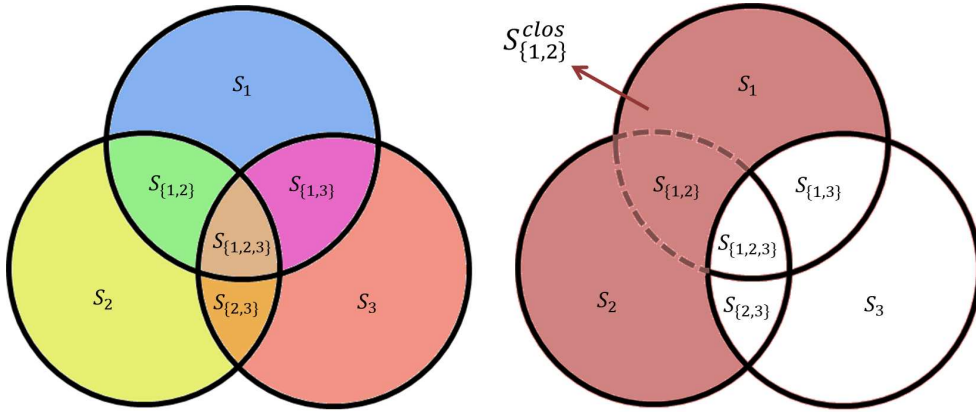


Figure A.11: Sobol' indices (left) and closed Sobol' indices (right).

## Appendix B. Axioms of Shapley values

Consider a game with  $d$  players, and let  $\text{val}(A) \in \mathbb{R}$  be the cost function quantifying the production value of a coalition (i.e., set of players)  $A \in \mathcal{P}_d$ , under the assumption that  $\text{val}(\emptyset) = 0$ . The Shapley value  $\phi_j = \phi_j(\text{val}), j = 1, \dots, d$  attributed to each player can be defined by the following set of axioms:

1. (Efficiency)  $\sum_{j=1}^d \phi_j = \text{val}(\{1, \dots, d\})$ , meaning that the sum of the allocated values have to be equal to the value produced by the cooperation of all the players.

2. (Symmetry) If  $\text{val}(A \cup \{i\}) = \text{val}(A \cup \{j\})$  for all  $A \in \mathcal{P}_d$ , then  $\phi_i = \phi_j$ , meaning that if two players allow for the same contribution to every coalition, their attribution should be the same.
3. (Dummy) If  $\text{val}(A \cup \{i\}) = \text{val}(A)$  for all  $A \in \mathcal{P}_d$ , then  $\phi_i = 0$ , meaning that if a player does not contribute to the production of resources for all coalition, he should not be attributed any resources.
4. (Additivity) If  $\text{val}$  and  $\text{val}'$  have Shapley Values  $\phi$  and  $\phi'$  respectively, then the game with cost function  $\text{val} + \text{val}'$  has Shapley values  $\phi_j + \phi'_j$  for  $j \in \{1, \dots, d\}$ .

These four axioms guarantee a cooperative allocation of  $\text{val}(\{1, \dots, d\})$ . The unique attribution method that satisfies these four axioms are the Shapley values [31], defined by:

$$\phi_j = \frac{1}{d} \sum_{A \subset -j} \binom{d-1}{|A|}^{-1} (\text{val}(A \cup \{j\}) - \text{val}(A)), \quad j = 1, \dots, d \quad (\text{B.1})$$

where  $\{-j\} = \{1, \dots, d\} \setminus \{j\}$ . One can additionally remark that  $\phi_j(\text{val})$  is a linear operator, meaning that for some constant  $c \in \mathbb{R}$ ,  $\phi_j(c \times \text{val}) = c \times \phi_j(\text{val})$ .

## Appendix C. Mathematical proofs

### Appendix C.1. Positivity of the $(\ell^1)$ -target Shapley effects

Let  $A \subseteq \{1, \dots, d\} \setminus \{j\}$ , for  $j \in \{1, \dots, d\}$ . In order to show that the  $(\ell^1)$ -target Shapley effects are positive, one needs to prove that:

$$\text{T-S}_{A \cup \{j\}}^{\ell^1} \geq \text{T-S}_A^{\ell^1}. \quad (\text{C.1})$$

In [49], it was shown that the following property holds:

$$\eta_{A \cup \{j\}} \geq \eta_A \quad (\text{C.2})$$

with  $\eta_A$  being defined in Eq. (18). From the definition of  $\text{T-S}_A^{\ell^1}$ ,

$$\text{T-S}_A^{\ell^1} = \frac{2}{\mathbb{E} \left[ \left| \mathbb{1}_{\mathcal{F}_i}(X) - \mathbb{E}[\mathbb{1}_{\mathcal{F}_i}(X)] \right| \right]} \eta_A, \quad (\text{C.3})$$

one gets immediately the property C.1.

### Appendix C.2. Positivity of the $(\ell^2)$ -target Shapley effects

Let  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ , be a real-valued random vector admitting a probability measure  $P_X$  on the usual real measurable space. Let  $L^2(P_X)$  be the functional space such that, for a measurable function  $f$ ,  $\|f\|_{L^2} \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} f^2(x) dP_X(x) < +\infty$ . Let  $G(\cdot) \in L^2$  be the studied numerical model, and denote the random variable  $Y = G(X)$  be the model output (or  $Y = \mathbb{1}_{G(X) > t}(X)$  the TSA variable of interest, without loss of generality). Let  $A \subseteq \{1, \dots, d\} \setminus \{j\}$  be the indices of the subset of inputs  $X_A$  and  $j \in \{1, \dots, d\}$ . In order to show that  $\text{T-Sh}_j \geq 0$ , one needs to prove that:

$$\text{T-S}_{A \cup \{j\}} - \text{T-S}_A \geq 0 \quad (\text{C.4})$$

which is equivalent to

$$\mathbb{V}\left(\mathbb{E}[Y|X_A]\right) \leq \mathbb{V}\left(\mathbb{E}[Y|X_{A \cup \{j\}}]\right). \quad (\text{C.5})$$

From the Pythagorean theorem, one has:

$$\|Y\|_{L^2} = \|\mathbb{E}[Y | X_A]\|_{L^2} + \|Y - \mathbb{E}[Y | X_A]\|_{L^2}, \quad (\text{C.6})$$

which is equivalent to

$$\mathbb{E}[Y^2] = \mathbb{E}\left[\left(\mathbb{E}[Y | X_A]\right)^2\right] + \mathbb{E}\left[\left(Y - \mathbb{E}[Y | X_A]\right)^2\right]. \quad (\text{C.7})$$

By removing  $(\mathbb{E}[Y])^2$  to both sides of the equality, one obtains:

$$\mathbb{V}\left(\mathbb{E}[Y|X_A]\right) = \mathbb{V}(Y) - \|Y - \mathbb{E}[Y|X_A]\|_{L^2}^2. \quad (\text{C.8})$$

By using the formula  $\mathbb{E}[Y | X_A] = \underset{Z \in \sigma(X_A)}{\operatorname{argmin}} \|Y - Z\|_{L^2}$ , with  $\sigma(X_A)$  being the span of  $X_A$ , we deduce that  $\mathbb{E}[Y | X_A] \leq \mathbb{E}[Y | X_{A \cup \{j\}}]$  since  $\sigma(X_A) \subseteq \sigma(X_{A \cup \{j\}})$ . This leads to

$$\mathbb{V}(Y) - \|Y - \mathbb{E}[Y | X_A]\|_{L^2}^2 \leq \mathbb{V}(Y) - \|Y - \mathbb{E}[Y | X_{A \cup \{j\}}]\|_{L^2}^2. \quad (\text{C.9})$$

Finally, from Eq. (C.8) and Eq. (C.9), we obtain

$$\mathbb{V}\left(\mathbb{E}[Y | X_A]\right) \leq \mathbb{V}\left(\mathbb{E}[Y | X_{A \cup \{j\}}]\right) \quad (\text{C.10})$$

which concludes the proof.

## Appendix D. Minimal R code examples for the estimation methods

### Appendix D.1. Monte Carlo sampling estimator

```
#Packages
library(sensitivity)
library(mvtnorm)
library(condMVNorm)

#Model definition
model.linear <- function(X) as.numeric(apply(X,1,sum)>0)

#Parameters
d <- 3
mu <- rep(0,d)
sig <- c(1,1,2)
ro <- 0.9
Cormat <- matrix(c(1,0,0,0,1,ro,0,ro,1),d,d)
Covmat <- ( sig %*% t(sig) ) * Cormat
```

```

#Total and marginal simulation function
Xall <- function(n) mvtnorm::rmvnorm(n,mu,Covmat)

#Conditional simulation function
Xset <- function(n, Sj, Sjc, xjc){
  if (is.null(Sjc)){
    if (length(Sj) == 1){ rnorm(n,mu[Sj],sqrt(Covmat[Sj,Sj]))
    }else{
      mvtnorm::rmvnorm(n,mu[Sj],Covmat[Sj,Sj])
    }
  }else{
    condMVNorm::rcmvnorm(n,
                        mu,
                        Covmat,
                        dependent.ind=Sj,
                        given.ind=Sjc,
                        X.given=xjc)
  }
}

#(l2)-target Shapley effects estimation
l2_tse.mc <- shapleyPermEx(model = modlin,
                          Xall=Xall,
                          Xset=Xset,
                          d=d,
                          Nv=1e4,
                          No = 1e3,
                          Ni = 3)

#Plot the results
print(l2_tse.mc)

#(l2)-target Shapley effects estimation with random permutations
l2_tse.mc.randperm<-shapleyPermRand(model = modlin,
                                    Xall=Xall,
                                    Xset=Xset,
                                    d=d,
                                    Nv=1e4,
                                    No = 1e3,
                                    Ni = 3,
                                    m=5)

#Plot the results
plot(l2_tse.mc.randperm)

```

Listing 1: Minimal R code example for the Monte Carlo estimation.

*Appendix D.2. Nearest-neighbor estimator*

```
#Packages
library(sensitivity)
library(mvtnorm)

#Random sample of inputs-output
X<-rmvnorm(2000, rep(0,3), diag(3))
Y<-rbinom(2000, 1, 0.7)

#(l2)-target Shapley effects estimation
l2_tse.knn<-sobolshap_knn(model=NULL,
                        X=X)
tell(l2_tse.knn, Y)

#Plot the results
plot(l2_tse.knn)

#(l2)-target Shapley effects estimation with random permutations
l2_tse.knn.randperm<-sobolshap_knn(model=NULL,
                                   X=X,
                                   rand.perm=T,
                                   n.perm=5)
tell(l2_tse.knn.randperm, Y)

#Plot the results
plot(l2_tse.knn.randperm)
```

Listing 2: Minimal R code example for the nearest-neighbor estimation.