

# Pratique de l'analyse de sensibilité : comment évaluer l'impact des entrées aléatoires sur la sortie d'un modèle mathématique\*

Julien JACQUES<sup>a</sup>

## Abstract

Global sensitivity analysis (SA) analyze a mathematical model by studying the impact of the variability of model input factors on the output variable. Determining the inputs responsible for this variability with sensitivity indices, SA allows to reduce the variance of output if it is synonymous with vagueness, or ease the model by setting the entries whose variability does not influence the output variable. We present in this paper the main sensitivity indices, based on the assumption of independent input variables, their estimates, and then address the case of non-independent input models. Two numerical applications illustrate the interpretation of sensitivity indices in the case of models with independent and dependent inputs.

## Résumé

L'analyse de sensibilité globale (AS) permet d'analyser un modèle mathématique en étudiant l'impact de la variabilité des facteurs d'entrée du modèle sur la variable de sortie. Déterminant les entrées responsables de cette variabilité à l'aide d'indices de sensibilité, l'AS permet de prendre les mesures nécessaires pour diminuer la variance de la sortie si celle-ci est synonyme d'imprécision, ou encore d'alléger le modèle en fixant les entrées dont la variabilité n'influe pas la variable de sortie. Nous présentons dans ce document les principaux indices de sensibilité, basés sur l'hypothèse d'indépendance des variables d'entrée, leurs estimations, puis abordons le cas des modèles à entrées non indépendantes. Deux applications numériques illustrent l'interprétation des indices de sensibilité dans le cas de modèle à entrées indépendantes et dépendantes.

*MSC 2009 subject classifications.* 49Q12.

*Mots clés.* Analyse de sensibilité, décomposition de la variance, indice de Sobol, entrée dépendante.

---

\*Preprint.

<sup>a</sup>Laboratoire P. Painlevé, UMR 8524 CNRS Université Lille I, Bât M2, Cité Scientifique, F-59655 Villeneuve d'Ascq Cedex, France.

## 1 Introduction : les objectifs de l'analyse de sensibilité

Considérons un modèle mathématique qui, à un ensemble de variable d'entrée aléatoire  $\mathbf{X}$ , fait correspondre, via une fonction  $f$  déterministe, une variable de sortie  $Y$  (ou réponse) aléatoire :

$$\begin{aligned} f : \mathbb{R}^p &\rightarrow \mathbb{R} \\ \mathbf{X} &\mapsto Y = f(\mathbf{X}) \end{aligned} \quad (1)$$

La fonction  $f$  du modèle peut être très complexe (système d'équation différentielle ...), et est en pratique évaluée à l'aide d'un code informatique, plus ou moins onéreux en temps de calcul. L'ensemble de variable d'entrée  $\mathbf{X} = (X_1, \dots, X_p)$  regroupe toute l'entité considérée comme aléatoire dans le modèle.

L'analyse de sensibilité étudie comment des perturbations sur les variables d'entrée du modèle engendrent des perturbations sur la variable réponse. L'auteur intéressé par un ouvrage de référence pourra se référer à [11]. Il est possible de grouper les méthodes d'analyse de sensibilité en trois classes : la méthode de *screening*, qui consiste en une analyse qualitative de la sensibilité de la variable de sortie aux variables d'entrée, la méthode d'analyse locale [17], qui évaluent quantitativement l'impact d'une petite variation autour d'une valeur donnée de l'entrée, et enfin la méthode d'analyse de sensibilité globale, qui s'intéressent à la variabilité de la sortie du modèle dans l'intégralité de son domaine de variation. L'analyse de sensibilité globale étudie comment la variabilité de l'entrée se répercute sur celle de la sortie, en déterminant quelle part de variance de la sortie est due à telle entrée ou tel ensemble d'entrée. Si l'analyse de sensibilité locale s'intéresse plus à la valeur de la variable réponse, l'analyse de sensibilité globale s'intéresse quant à elle à sa variabilité. Nous nous intéressons dans ce document à l'analyse de sensibilité globale et omettrons donc par la suite l'adjectif global.

Les enjeux de l'analyse de sensibilité peuvent être multiples : validation d'une méthode ou d'un code de calcul, orientation de l'effort de recherche et développement, ou encore justification en terme de coût d'un dimensionnement ou d'une modification d'un système. Nous décrivons ci-après les principales questions auxquelles l'analyse de sensibilité permet d'apporter des éléments de réponse.

**Les ambitions de l'analyse de sensibilité** Au cours de l'élaboration, de la construction ou de l'utilisation d'un modèle mathématique, l'analyse de sensibilité peut s'avérer être un outil précieux. En effet, en étudiant comment la réponse du modèle réagit aux variations de ses variables d'entrée, l'analyse de sensibilité permet de répondre à un certain nombre de questions.

1. Le modèle est-il bien fidèle au phénomène/processus modélisé ?  
En effet, si l'analyse exhibe une forte influence d'une variable d'entrée

habituellement connue comme non influente, il sera nécessaire de remettre en cause la qualité du modèle ou (et) la véracité de nos connaissances sur l'impact réel de la variable d'entrée.

2. Quelles sont les variables qui contribuent le plus à la variabilité de la réponse du modèle ?

Si cette variabilité est synonyme d'imprécision sur la valeur prédite de la sortie, il sera alors possible d'améliorer la qualité de la réponse du modèle à moindre coût. En effet, la variabilité de la sortie du modèle pourra être diminuée en concentrant l'effort sur la réduction de la variabilité d'entrée la plus influente. Il doit être précisé que cela n'est pas toujours possible, notamment lorsque la variabilité d'une variable d'entrée est intrinsèque à la nature de la variable et non due à un manque d'information ou à de l'imprécision de mesure.

3. Quelles sont au contraire les variables les moins influentes ?

Il sera possible de les considérer comme des paramètres déterministes, en les fixant par exemple à leur espérance, et obtenir ainsi un modèle plus léger avec moins de variables d'entrée. Dans le cas d'un code informatique, il sera possible de supprimer la partie de code qui n'ont aucune influence sur la valeur et la variabilité de la réponse.

4. Quelles variables, ou quel groupe de variables, interagissent avec quelle (quelles) autre(s) ?

L'analyse de sensibilité peut permettre de mieux appréhender et comprendre le phénomène modélisé, en éclairant les relations entre les variables d'entrée.

Bon nombre de publications sur le sujet explicitent et illustrent cet objectif. On pourra se référer notamment aux travaux de Saltelli et al. [12, 13, 15].

La section suivante présente l'indice de sensibilité défini pour un modèle à variable d'entrée indépendante, ainsi que leur méthode d'estimation. La section 3 s'intéresse aux modèles à entrée non indépendante, et présente deux types d'indices utilisables dans ce cas. Enfin, la section 4 présente deux applications simulées illustrant l'intérêt et l'interprétation de l'indice de sensibilité, dans le cas de modèles à entrée indépendante et non indépendante.

## 2 Indicateurs de sensibilité pour modèles à entrées indépendantes

Nous supposons dans cette section que les variables d'entrée  $\mathbf{X} = (X_1, \dots, X_p)$  du modèle sont indépendantes.

## 2.1 Préambule : cas du modèle linéaire

Supposons que le modèle étudié soit linéaire, et qu'il s'écrive sous la forme suivante :

$$Y = \mu_0 + \sum_{i=1}^p \beta_i X_i. \quad (2)$$

Comme les variables  $X_i$  sont supposées indépendantes, la variance de  $Y$  s'écrit alors :

$$V(Y) = \sum_{i=1}^p \beta_i^2 V(X_i);$$

où  $\beta_i^2 V(X_i)$  est la part de variance due à la variable  $X_i$ . La sensibilité de  $Y$  à  $X_i$  peut donc simplement être quantifiée par le rapport de la part de variance due à  $X_i$  sur la variance totale. On définit ainsi l'indice de sensibilité *SRC* (*Standardized Regression Coefficient*) :

$$SRC_i = \frac{\beta_i^2 V(X_i)}{V(Y)}. \quad (3)$$

Il exprime la part de variance de la réponse  $Y$  due à la variance de la variable  $X_i$ . Cet indice *SRC*, toujours positif ( $SRC \in [0; 1]$ ), est en outre le carré du coefficient de corrélation linéaire entre la réponse du modèle et la variable d'entrée.

## 2.2 Les indices de Sobol

Plaçons nous désormais dans le cas d'une fonction  $f$  dont la forme analytique n'est pas connue. Pour apprécier l'importance d'une variable d'entrée  $X_i$  sur la variance de la sortie  $Y$ , nous étudions à combien la variance de  $Y$  décroît si on fixe la variable  $X_i$  à une valeur  $x_i$  :  $V(Y|X_i = x_i)$ . Le problème de cet indicateur est le choix de la valeur  $x_i$  de  $X_i$ , que l'on résout en considérant l'importance de cette quantité pour toute la valeur possible de  $x_i$  :  $\mathbf{E}[V(Y|X_i)]$ . Ainsi, plus la variable  $X_i$  sera importante vis-à-vis de la variance de  $Y$ , plus cette quantité sera petite. Étant donné la formule de la variance totale  $V(Y) = V(\mathbf{E}[Y|X_i]) + \mathbf{E}[V(Y|X_i)]$ , nous pouvons utiliser de façon équivalente la quantité

$$V(\mathbf{E}[Y|X_i]);$$

qui sera d'autant plus grande que la variable  $X_i$  sera importante vis-à-vis de la variance de  $Y$ . Afin d'utiliser un indicateur normalisé, nous définissons l'indice de sensibilité de  $Y$  à  $X_i$  :

$$S_i = \frac{V(\mathbf{E}[Y|X_i])}{V(Y)}. \quad (4)$$

Cet indice est appelé **indice de sensibilité de premier ordre** par Sobol [16], *correlation ratio* par McKay [5], ou encore *importance measure*. Il quantifie la

en ibilité de la ortie  $Y$  à la variable d'entrée  $X_i$ , ou encore la part de variance de  $Y$  due à la variable  $X_i$ .

**Remarque.** Dans le cas du modèle linéaire (2), cet indice de en ibilité est égal à l'indice  $SRC$ , puisque  $V(\mathbf{E}[Y|X_i]) = V(\mathbf{E}[X_i]) = \mathbf{E}[X_i^2] - (\mathbf{E}[X_i])^2 = V(X_i)$ .

Sobol [16] a introduit cet indice de en ibilité en décomposant la fonction  $f$  du modèle en somme de fonction de dimension croissante :

$$Y = f(X_1, \dots, X_p) \quad (5)$$

$$= f_0 + \sum_{i=1}^p f_i(X_i) + \sum_{1 \leq i < j \leq p} f_{ij}(X_i, X_j) + \dots + f_{1\dots p}(X_1, \dots, X_p) \quad (6)$$

où

$$\begin{aligned} f_0 &= \mathbf{E}[Y]; \\ f_i(X_i) &= \mathbf{E}[Y|X_i] - \mathbf{E}[Y]; \\ f_{i,j}(X_i, X_j) &= \mathbf{E}[Y|X_i, X_j] - \mathbf{E}[Y|X_i] - \mathbf{E}[Y|X_j] + \mathbf{E}[Y]; \\ f_{i,j,k}(X_i, X_j, X_k) &= \mathbf{E}[Y|X_i, X_j, X_k] - \mathbf{E}[Y|X_i, X_j] - \mathbf{E}[Y|X_i, X_k] \\ &\quad - \mathbf{E}[Y|X_j, X_k] + \dots \end{aligned}$$

La variance de  $Y$ ,  $V$ , peut alors se décomposer selon le théorème suivant.

**Théorème.** Décomposition de Sobol de la variance.

La variance du modèle à entrée indépendante (1) se décompose en :

$$V = \sum_{i=1}^p V_i + \sum_{1 \leq i < j \leq p} V_{ij} + \dots + V_{1\dots p} \quad (7)$$

où

$$\begin{aligned} V_i &= V(\mathbf{E}[Y|X_i]); \\ V_{ij} &= V(\mathbf{E}[Y|X_i, X_j]) - V_i - V_j; \\ V_{ijk} &= V(\mathbf{E}[Y|X_i, X_j, X_k]) - V_{ij} - V_{ik} - V_{jk} - V_i - V_j - V_k; \\ &\dots \\ V_{1\dots p} &= V - \sum_{i=1}^p V_i - \sum_{1 \leq i < j \leq p} V_{ij} - \dots - \sum_{1 \leq i_1 < \dots < i_{p-1} \leq p} V_{i_1 \dots i_{p-1}} \end{aligned}$$

expriment la sensibilité de la variance de  $Y$  aux variables  $X_i, X_j$  et  $X_k$  qui n'est pas prise en compte dans l'effet de variable seule et de interaction deux à deux. Et ainsi de suite jusqu'à l'ordre  $p$ .

**L'interprétation de ces indices est facile**, puisque grâce à (7), **leur somme est égale à 1**, et étant toujours positif, plus l'indice sera grand (proche de 1), plus la variable aura d'importance.

Le nombre d'indices de sensibilité ainsi construit, de l'ordre 1 à l'ordre  $p$ , est égal à  $2^p - 1$ . Lorsque le nombre de variables d'entrée  $p$  est trop important, le nombre d'indices de sensibilité explose. L'estimation et l'interprétation de tous ces indices deviennent vite impossibles. Homma et Saltelli [2] ont alors introduit des indices de sensibilité totaux, qui expriment la sensibilité totale de la variance  $Y$  à une variable, c'est-à-dire la sensibilité à cette variable sous toute forme (sensibilité à la variable seule et sensibilité aux interactions de cette variable avec d'autres variables).

**L'indice de sensibilité total**  $S_{T_i}$  à la variable  $X_i$  est défini comme la somme de tous les indices de sensibilité relatifs à la variable  $X_i$  :

$$S_{T_i} = \sum_{k \neq i} S_k \quad (8)$$

où  $\#i$  représente tous les ensembles d'indices contenant l'indice  $i$ .

Exemple : pour un modèle à trois variables d'entrée  $S_{T_1} = S_1 + S_{12} + S_{13} + S_{123}$ .

### 2.3 Estimation des indices de Sobol

**Estimation de Monte Carlo** Dans beaucoup de problèmes scientifiques, on est amené à calculer une intégrale du type

$$I = \int_D f(\mathbf{x}) d\mathbf{x};$$

où  $D$  est un espace de plus ou moins grande dimension, et  $f$  une fonction (intégrable). Soit  $x_1, \dots, x_N$  la réalisation d'un  $N$ -échantillon d'une variable aléatoire uniforme sur  $D$ . Nous supposons cet échantillon pris de manière totalement aléatoire (échantillonnage aléatoire). Une approximation de  $I$  par la méthode de Monte Carlo est faite par :

$$\hat{I}_N = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i);$$

La convergence (presque sûre) de  $\hat{I}_N$  vers  $I$  découle directement de la loi forte des grands nombres. Cette méthode d'estimation permet alors d'estimer l'espérance de toute fonction d'une variable aléatoire de densité quelconque par

$$\hat{\mathbf{E}}[f(X)] = \frac{1}{N} \sum_{i=1}^N f(x_i);$$

où  $(X_i)_{i=1..N}$  est un  $N$ -échantillon de réalisation de la variable aléatoire  $X$ . Le taux de convergence d'une méthode de Monte Carlo est en  $\mathcal{O}(N^{-\frac{1}{2}})$ .

Bon nombre de méthodes alternatives ont été proposées pour améliorer la convergence, parmi lesquelles la méthode de simulation pseudo-probabiliste<sup>1</sup>, comme l'échantillonnage stratifié ou par hypercube latin (*LHS*) [6], la méthode de Quasi-Monte Carlo [7], ou encore la méthode de Quasi-Monte Carlo Randomisé [9]. L'échantillonnage stratifié consiste à découper l'espace de variable d'entrée en petites cellules disjointes, puis à échantillonner au sein de chacune de ces petites cellules. L'échantillonnage *LHS* est basé sur le même principe, en assurant que le découpage a défini des petites cellules équiprobables, et que chaque petite cellule est bien échantillonnée; le quadrillage se fait dans le cube unité, pour un tirage aléatoire d'échantillon uniforme, puis ce échantillon est transformé via la fonction de répartition inverse. Les méthodes de Quasi-Monte Carlo ont de meilleures performances que la méthode de Monte Carlo. Ces méthodes définissent des séquences d'échantillon déterministes qui ont une densité plus faible que la séquence aléatoire, c'est-à-dire qu'elles ont une meilleure répartition uniforme dans l'espace de variable d'entrée. Ces méthodes de Quasi-Monte Carlo permettent d'obtenir une convergence plus rapide en  $\mathcal{O}(N^{-1}(\log N)^{p-1})$  (ou de condition relativement faible de régularité de  $f$ ). Les méthodes de Quasi-Monte Carlo Randomisé, ou certaines conditions peu restrictives sur  $f$ , ont un taux de convergence en  $\mathcal{O}(N^{-\frac{3}{2}}(\log N)^{\frac{p-1}{2}})$ , et permettent une approximation de l'erreur d'estimation. Owen [8] présente ces méthodes comme une ré-allocation de séquence utilisée dans la méthode de Quasi-Monte Carlo : on prend la séquence déterministe  $a_i$  de cette dernière, et on la transforme en variable aléatoire  $X_i$ . Cette transformation se fait par exemple par  $X_i = a_i + U \bmod 1$ ; où  $U \sim U[0;1]^p$ .

### Estimation des indices de sensibilité par Monte Carlo

Considérons un  $N$ -échantillon  $\tilde{X}_{(N)} = (X_{k1} : \dots : X_{kp})_{k=1..N}$  de réalisation de variable d'entrée  $(X_1 : \dots : X_p)$ . L'espérance de  $Y$ ,  $\mathbf{E}[Y] = f_0$ , et la variance,  $V(Y) = V$ , sont estimées par :

$$\hat{f}_0 = \frac{1}{N} \sum_{k=1}^N f(X_{k1} : \dots : X_{kp}); \quad \text{et} \quad \hat{V} = \frac{1}{N} \sum_{k=1}^N f^2(X_{k1} : \dots : X_{kp}) - \hat{f}_0^2 \quad (9)$$

L'estimation de l'indice de sensibilité nécessite l'estimation de la variance conditionnelle. Nous présentons une technique d'estimation due à Sobol [16].

L'estimation de l'indice de sensibilité de premier ordre (4) consiste à estimer la quantité :

$$V_i = V(\mathbf{E}[Y|X_i]) = \underbrace{\mathbf{E}[\mathbf{E}[Y|X_i]^2]}_{U_i} - \mathbf{E}[\mathbf{E}[Y|X_i]]^2 = U_i - \mathbf{E}[Y]^2;$$

<sup>1</sup> «pseudo» puisqu'elle consiste en un échantillonnage non totalement aléatoire

la variance de  $Y$  étant estimée classiquement par (9). Sobol propose d'estimer la quantité  $U_i$ , c'est-à-dire l'espérance du carré de l'espérance de  $Y$  conditionnellement à  $X_i$ , comme une espérance classique, mais en tenant compte du conditionnement à  $X_i$  en faisant varier entre les deux appels à la fonction  $f$  toute la variable sauf la variable  $X_i$ . Ceci nécessite deux échantillons de réalisation de variable d'entrée, que nous notons  $\tilde{X}_{(N)}^{(1)}$  et  $\tilde{X}_{(N)}^{(2)}$  :

$$\hat{U}_i = \frac{1}{N} \sum_{k=1}^N f \left( x_{k1}^{(1)}; \dots; x_{k(i-1)}^{(1)}; x_{ki}^{(1)}; x_{k(i+1)}^{(1)}; \dots; x_{kp}^{(1)} \right) \\ f \left( x_{k1}^{(2)}; \dots; x_{k(i-1)}^{(2)}; x_{ki}^{(1)}; x_{k(i+1)}^{(2)}; \dots; x_{kp}^{(2)} \right) :$$

L'indice de sensibilité de premier ordre est alors estimé par :

$$\hat{S}_i = \frac{\hat{V}_i}{\hat{V}} = \frac{\hat{U}_i - \hat{f}_0^2}{\hat{V}} :$$

Pour l'indice de sensibilité de second ordre  $S_{ij} = \frac{V_{ij}}{V}$ , où :

$$V_{ij} = \mathbf{V}(\mathbf{E}[Y|X_i; X_j]) - V_i - V_j = U_{ij} - \mathbf{E}[Y]^2 - V_i - V_j ;$$

nous estimons la quantité  $U_{ij} = \mathbf{E}[\mathbf{E}[Y|X_i; X_j]^2]$  de la même manière, en faisant varier entre les deux appels à la fonction toute la variable sauf  $X_i$  et  $X_j$  :

$$\hat{U}_{ij} = \frac{1}{N} \sum_{k=1}^N f \left( x_{k1}^{(1)}; \dots; x_{k(i-1)}^{(1)}; x_{ki}^{(1)}; x_{k(i+1)}^{(1)}; \dots; x_{k(j-1)}^{(1)}; x_{kj}^{(1)}; x_{k(j+1)}^{(1)}; \dots; x_{kp}^{(1)} \right) \\ \times f \left( x_{k1}^{(2)}; \dots; x_{k(i-1)}^{(2)}; x_{ki}^{(1)}; x_{k(i+1)}^{(2)}; \dots; x_{k(j-1)}^{(2)}; x_{kj}^{(1)}; x_{k(j+1)}^{(2)}; \dots; x_{kp}^{(2)} \right) :$$

L'indice  $S_{ij}$  est alors estimé par :

$$\hat{S}_{ij} = \frac{\hat{U}_{ij} - \hat{f}_0^2 - \hat{V}_i - \hat{V}_j}{\hat{V}} :$$

Et ainsi de suite pour l'indice de sensibilité d'ordre supérieur.

**Remarque.** L'estimation de l'indice de sensibilité d'ordre  $i$ , ( $1 < i \leq \rho$ ), nécessite l'estimation de l'indice de sensibilité d'ordre 1 à  $i-1$ .

Par contre, les indices de sensibilité totaux peuvent être estimés directement. En effet, on remarque facilement que l'indice de sensibilité total peut s'écrire

$$S_{Ti} = 1 - \frac{\mathbf{V}(\mathbf{E}[Y|X_{-i}])}{\mathbf{V}(Y)} = 1 - \frac{V_{-i}}{V} :$$

où  $V_{-i}$  est la variance de l'espérance de  $Y$  conditionnellement à toute la variable sauf  $X_i$ .  $V_{-i}$  est alors estimée comme  $V_i$ , sauf qu'au lieu de faire varier



toute le variable sur  $X_i$ , nous ne faisons varier uniquement  $X_i$ . Ainsi, pour estimer  $V_i = \mathbf{E}[\mathbf{E}[Y|X_i]^2] - \mathbf{E}[\mathbf{E}[Y|X_i]]^2 = U_i - \mathbf{E}[Y]^2$ , on estime  $U_i$  par :

$$\hat{U}_i = \frac{1}{N} \sum_{k=1}^N \left( f\left(x_{k1}^{(1)}; \dots; x_{k(i-1)}^{(1)}; x_{ki}^{(1)}; x_{k(i+1)}^{(1)}; \dots; x_{kp}^{(1)}\right) \right. \\ \left. f\left(x_{k1}^{(1)}; \dots; x_{k(i-1)}^{(1)}; x_{ki}^{(2)}; x_{k(i+1)}^{(1)}; \dots; x_{kp}^{(1)}\right); \right.$$

et on obtient

$$\hat{S}_{T_i} = 1 - \frac{\hat{U}_i - \hat{f}_0^2}{\hat{V}}$$

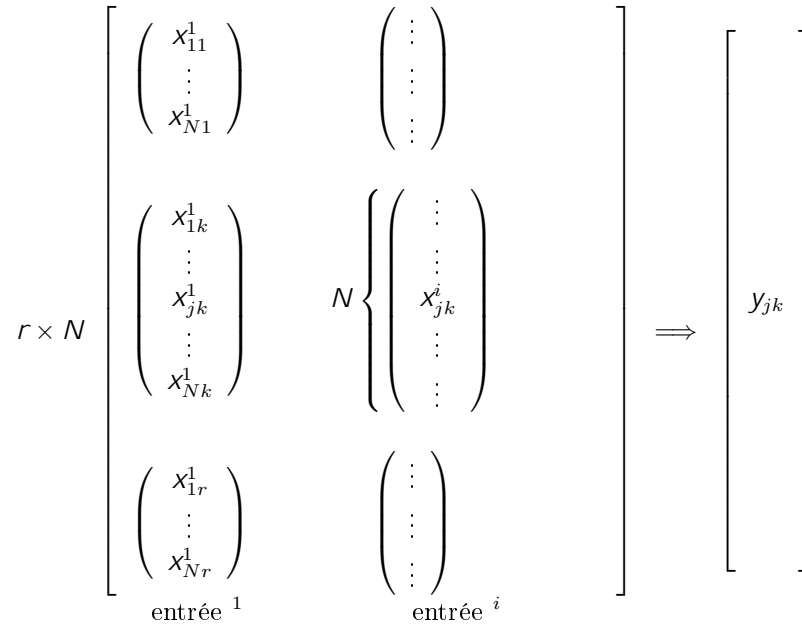
**Quels indices estimer : stratégie à adopter** En utilisant une taille d'échantillon de Monte Carlo de  $N$ , le nombre réel de simulation de variable d'entrée nécessaire à l'estimation de l'indice de variabilité est  $2N$ , puisque cette estimation nécessite deux jeux de simulation. Le nombre d'appel à la fonction du modèle est alors  $N \times (k + 1)$ , où  $k$  est le nombre d'indice estimé. Pour un modèle à  $p$  variable d'entrée, l'estimation de tous les indices de variabilité nécessite  $N \times (2^p)$  appel à la fonction. En revanche, n'estimer que le indice de premier ordre et le indice totaux ne demande que  $N \times (2p + 1)$  appel. Il conviendra donc d'estimer dans un premier temps le indice de premier ordre et le indice totaux. S'il existe de l'écart important entre ces deux indices, c'est que la part de l'interaction est non négligeable et il peut être utile d'estimer le indice d'ordre intermédiaire. Dans le cas contraire, l'effet de variable d'entrée sera principalement de premier ordre et il ne sera pas utile d'intéresser aux indices d'ordre intermédiaire.

En pratique, une taille d'échantillon de l'ordre de 10000 sera suffisante pour estimer le indice de variabilité d'un modèle comportant une dizaine de variable d'entrée. En outre, il sera possible d'estimer la variabilité de l'estimateur obtenu par bootstrap. Lorsque le modèle demande un temps d'exécution important, il est illusoire de vouloir utiliser de telle taille d'échantillon en un temps raisonnable. On a en général recouru à une approximation de la fonction  $f$  (surface de réponse), permettant de faire de simulation intensive et donc d'estimer le indice de variabilité. Le lecteur intéressé par une revue de méthode de surface de réponse pour l'analyse de variabilité pourra se référer à [3] par exemple.

### 2.3.1 La méthode de McKay

La méthode d'estimation de l'indice de variabilité de premier ordre proposée par McKay, [5], est basée sur l'échantillonnage par hypercube latin répliqué (*r-LHSampling*). À partir d'un  $N$ -échantillon créé selon le plan d'échantillonnage par hypercube latin ( $N$  première ligne de la matrice ci-dessus), on crée  $r$  répliques (paquet de  $N$  lignes) en permutant indépendamment et aléatoirement les  $N$  valeurs de chaque variable (i.e. colonne). La réunion de ces  $r$  répliques

donnera  $N \times r$  échantillon pour chaque variable. Ce schéma d'échantillonnage par hypercube latin répliqué peut être représenté par la figure 1.



$1 \leq j \leq N$  : N valeur de entrée (prise dans un intervalle équiprobable),  
 $1 \leq k \leq r$  : r permutation de N-vecteur de simulation de entrée,  
 $1 \leq i \leq p$  : p paramètre.

FIGURE 1 – Échantillonnage par hypercube latin répliqué.

Le moyenne suivante sont alors définies :

$$\bar{y}_{j.} = \frac{1}{r} \sum_{k=1}^r y_{jk} \quad \bar{y} = \frac{1}{N} \sum_{j=1}^N \bar{y}_{j.}$$

où  $\bar{y}_{j.}$  est la moyenne *inter* réplication et  $\bar{y}$  est la moyenne sur toute la valeur de  $y$ .

L'estimation de l'indice de sensibilité de premier ordre de la variable  $X_i$ , défini par (4) nécessite l'estimation de la quantité  $V(\mathbf{E}[Y|X_i])$  et  $V(Y)$ . La variance

totale  $V(Y)$  peut être estimée par :

$$\widehat{V^{(\cdot)}}(Y) = \frac{1}{r} \sum_{k=1}^r \underbrace{\frac{1}{N} \sum_{j=1}^N (y_{jk} - \bar{y}_{.k})^2}_{\widehat{V}_k(Y)}; \quad (10)$$

où  $\bar{y}_{.k} = \frac{1}{N} \sum_{j=1}^N y_{jk}$  et  $\widehat{V}_k(Y)$  ont le même dénominateur de la moyenne et de la variance de  $Y$  au sein de la réplication  $k$  (*intra* réplication). En utilisant la formule classique de l'analyse de la variance, pour une somme de carrés *intra* et *inter* réplication, qui s'écrit :

$$\begin{aligned} \sum_{k=1}^r \sum_{j=1}^N (y_{jk} - \bar{y})^2 &= \underbrace{\sum_{k=1}^r \sum_{j=1}^N (\bar{y}_{.k} - \bar{y})^2}_{inter} + \underbrace{\sum_{k=1}^r \sum_{j=1}^N (y_{jk} - \bar{y}_{.k})^2}_{intra} \\ &= N \sum_{k=1}^r (\bar{y}_{.k} - \bar{y})^2 + \sum_{k=1}^r \sum_{j=1}^N (y_{jk} - \bar{y}_{.k})^2; \end{aligned}$$

on a :

$$\widehat{V^{(\cdot)}}(Y) = \frac{1}{Nr} \sum_{k=1}^r \sum_{j=1}^N (y_{jk} - \bar{y})^2 - \frac{1}{r} \sum_{k=1}^r (\bar{y}_{.k} - \bar{y})^2;$$

Or, pour un échantillonnage *LHS*, comme  $\mathbf{E}[(\bar{y}_{.k} - \bar{y})^2]$  est en  $\frac{1}{N}$ , le dernier terme de cette égalité peut être considéré comme négligeable pour une taille d'échantillon  $N$  suffisamment grande. McKay propose alors l'estimation de la variance totale suivante :

$$\widehat{V}(Y) = \frac{1}{Nr} \sum_{j=1}^N \sum_{k=1}^r (y_{jk} - \bar{y})^2;$$

Soient  $\bar{Y}_j$  et  $\bar{Y}$  le variable aléatoire dont  $\bar{y}_j$  et  $\bar{y}$  ont la réalisation sur notre matrice d'échantillonnage. Comme :

$$\begin{aligned} \mathbf{E}[(\bar{Y}_j - \bar{Y})^2] \simeq V(\bar{Y}_j) &= V(\mathbf{E}[\bar{Y}_j | X_i]) + \mathbf{E}[V(\bar{Y}_j | X_i)] \\ &= V(\mathbf{E}[Y | X_i]) + \frac{1}{r} \mathbf{E}[V(Y | X_i)]; \end{aligned}$$

le terme  $V(\mathbf{E}[Y | X_i])$  est estimé par :

$$\frac{1}{N} \sum_{j=1}^N (\bar{Y}_j^{(i)} - \bar{y})^2 - \frac{1}{r} \frac{1}{Nr} \sum_{j=1}^N \sum_{k=1}^r (y_{jk}^{(i)} - \bar{y}_j)^2;$$

où  $\frac{1}{Nr} \sum_{j=1}^N \sum_{k=1}^r (y_{jk}^{(i)} - \bar{y}_j)^2$  est l'estimateur de  $\mathbf{E}[V(Y|X_i)]$ , avec  $y_{jk}^{(i)}$  et  $\bar{y}_j^{(i)}$  obtenu en fixant, pour la variable  $X_i$ , la  $r$  réplique, ( $X_{jk}^i$  constant sur  $k$ , c'est-à-dire  $X_{j1}^i = X_{j2}^i = \dots = X_{jr}^i$  pour tout  $1 \leq j \leq N$ ). L'indice de dépendabilité de premier ordre de la variable  $X_i$ , défini par (4) est alors estimé par :

$$S_i = \frac{r \sum_{j=1}^N (\bar{y}_j^{(i)} - \bar{y})^2 - \frac{1}{r} \sum_{j=1}^N \sum_{k=1}^r (y_{jk}^{(i)} - \bar{y}_j)^2}{\sum_{j=1}^N \sum_{k=1}^r (y_{jk} - \bar{y})^2}.$$

### 3 Modèles à entrées dépendantes

L'hypothèse de l'indépendance de facteur d'entrée faite précédemment est nécessaire pour garantir l'interprétabilité de l'indice (un indice d'ordre un n'exprime plus la dépendabilité à une unique variable si cette dernière est corrélée avec d'autres) et la validité de leur méthode d'estimation par Monte-Carlo (le produit multidimensionnel est évalué comme le produit d'intégrales unidimensionnelles).

Nous présentons dans cette section la stratégie possible pour réaliser une analyse de dépendabilité sur un modèle à variable d'entrée non indépendante.

#### 3.1 Indices multidimensionnels

Lorsque toute la variable d'entrée n'est pas dépendante, mais qu'elle peut être regroupée en cluster de variable dépendante (le variable au sein d'un cluster dépendante mais la variable de différent cluster sont indépendantes), il est possible de considérer des indices de dépendabilité multidimensionnel [4] qui expriment la dépendabilité de la variance de  $Y$  à un cluster de facteur.

Si par exemple les deux variables  $X_i$  et  $X_j$  sont dépendantes, mais indépendante du reste de autres variables, la dépendabilité à la variable bidimensionnelle  $(X_i; X_j)$  sera exprimé par l'indice multidimensionnel

$$S_{\bar{r}_i, jg} = \frac{V(E[Y|X_i; X_j])}{V(Y)}.$$

Il est possible de définir des indices d'ordre supérieur exprimant la dépendabilité de  $Y$  à l'interaction entre cette variable bidimensionnelle  $(X_i; X_j)$  et d'autres variables uni ou multidimensionnelles. Le cluster de variable étant indépendant entre eux, l'interprétabilité (et en particulier la sommation de l'indice de tout ordre à 1) est conservée.

L'estimation de ce indice peut être faite par Monte Carlo avec une approche

similaire à celle utilisée pour estimer le indice de sensibilité de Sobol classique (cf. [4] pour plus de détail).

### 3.2 Utilisation des indices de Sobol d'ordre 1

Lorsque l'analyse de sensibilité est menée dans le but de savoir quelle variable ou quel groupe de variable qui, une fois fixé, conduit à la plus grande réduction de la variance de  $Y$ , Saltelli et Tarantola [14] expliquent que le indice de sensibilité d'ordre un ont toujours le indicateur à utiliser en présence de corrélation. En effet, en présence de corrélation l'indice d'ordre un  $S_i$  n'exprime plus uniquement la sensibilité à une variable  $X_i$  mais également une partie de sensibilité aux variables avec lesquelles elle est corrélée, fixer  $X_i$  conduit à également jouer sur la distribution de variables avec lesquelles elle est corrélée, et donc conduit à réduire d'autant plus la variance de la réponse du modèle.

Si l'estimation de Monte-Carlo de indice de premier ordre présentée précédemment (section 2.3) n'est plus valable en l'absence d'indépendance entre les variables d'entrée, la méthode de McKay (section 2.3.1) est toujours valable. Néanmoins, cette méthode d'estimation est très gourmande en nombre d'évaluation de la fonction  $f$  du modèle, ce qui peut être problématique lorsque l'évaluation de  $f$  est coûteuse en temps de calcul. Nous présentons ci-après une méthode d'estimation par polynôme locaux réduisant considérablement ce nombre d'évaluation [1].

**Estimation par polynômes locaux** La méthode d'estimation de indice de sensibilité d'ordre 1 de Da Veiga [1] consiste à estimer dans un premier temps l'espérance  $Y$  conditionnellement à chaque variable d'entrée  $X_i$ , puis dans un second temps à estimer la variance de cette espérance conditionnelle pour obtenir l'estimateur de l'indice de sensibilité. L'avantage principal de cette méthode est qu'elle ne fait appel qu'à un nombre réduit d'appels à la fonction, contrairement à la méthode de McKay précédente.

Notons  $m_i(x) = E[Y|X_i = x]$ . On approche  $m(x)$  localement par un polynôme

$$m_i(z) \simeq \sum_{j=0}^p c_j(z - x)^j \quad \forall z \in \mathcal{V}(x)$$

où  $\mathcal{V}(x)$  un voisinage de  $x$ , symbolisé par une fonction noyau  $K$  (de paramètre d'échelle  $h$ ) pondérant l'estimation par moindres carrés :

$$= \operatorname{argmin}_{\beta} \sum_{j=1}^n \left( Y^j - \sum_{j=0}^p c_j(X_i^j - x)^j \right)^2 K \left( \frac{X_i^j - x}{h} \right);$$

avec  $(X_i^j; Y^j)_{j=1, n}$  un échantillon de réalisation du couple  $(X_i; Y)$ . Utilisant un second échantillon  $(\tilde{X}_i^j)_{j=1, n'}$  de réalisation de la variable  $X_i$ , indépendant du

premier, on peut estimer classiquement la variance de  $m_i(x)$  par :

$$U_i = \frac{1}{n^\theta - 1} \sum_{j=1}^{n'} (m_i(\tilde{X}_i^j) - \bar{m}_i)^2$$

où  $\bar{m}_i = \sum_{j=1}^{n'} m_i(\tilde{X}_i^j) = n^\theta$ . Il suffit alors de diviser par l'estimation de la variance de  $Y$  pour obtenir une estimation de l'indice de sensibilité d'ordre un  $S_i$ .

## 4 Outil logiciel sous R et illustrations numériques

Dans cette section, après avoir précisé le package **R** permettant de réaliser de l'analyse de sensibilité, nous présentons deux analyses de sensibilité de modèle simulé, dans le cas d'entrée indépendante puis non indépendante.

### 4.1 Outil logiciel sous R

**Package sensitivity** Le package `sensitivity` [10] du logiciel **R**, disponible sur le site du CRAN<sup>2</sup> permet de calculer le indice de sensibilité de Sobol pré-enté dans ce document, lorsque les variables d'entrée sont indépendantes. La fonction `sobol` permet de calculer le indice de tout ordre, tandis que la fonction `sobol2002` permet d'estimer le indice de premier ordre et d'ordre total à partir d'un nombre d'échantillon plus réduit que la fonction `sobol`. Ces deux fonctions retournent des intervalles de confiance estimés par boot trap.

**Package sensitivity-dependent** Un package `sensitivity-dependent` pour le logiciel **R**, disponible sur le site de l'auteur<sup>3</sup>, permet de calculer le indice de sensibilité multidimensionnel (fonction `sobol_multi`) et le indice de sensibilité de premier ordre par la méthode de McKay (fonction `mckay`) en présence de variables d'entrée dépendantes. La fonction `sobol_multi` fournit en outre une estimation de la variabilité de l'estimation par boot trap.

### 4.2 Illustration de modèles à entrées indépendantes

Nous considérons trois modèles à entrée indépendante :

- le modèle linéaire  $Y = 9X_1 + 6X_2 + 3X_3 + X_4$  avec  $X_i \sim \mathcal{U}[0;1]$ ,
- le benchmark d'Ihigami [11] :  $Y = \sin(X_1) + 7 \sin^2(X_2) + \frac{X_3^4}{10} \sin(X_1)$  où  $X_i \sim \mathcal{U}[-1;1]$  pour  $i = 1;2;3$ ,
- le benchmark de Sobol [11] :  $Y = \prod_{j=1}^8 \frac{j^4 X[j] - 2j + a[j]}{1+a[j]}$  avec  $a = [0;1;4;5;9;99;99;99;99]$  et  $X_i \sim \mathcal{U}[0;1]$ .

<sup>2</sup><http://cran.r-project.org/web/packages/sensitivity/>

<sup>3</sup><http://labomath.univ-lille1.fr/~jacques/>

variable	indice ordre 1	inter. conf.	indice total	inter. conf.
modèle linéaire				
$X_1$	0.593	[0.523,0.661]	0.684	[0.603,0.753]
$X_2$	0.252	[0.200,0.298]	0.300	[0.250,0.345]
$X_3$	0.071	[0.055,0.097]	0.068	[0.045,0.089]
$X_4$	0.009	[0,0.019]	0.007	[0,0.016]
I higami benchmark				
$X_1$	0.305	[0.269,0.338]	0.578	[0.546,0.607]
$X_2$	0.4356	[0.403,0.462]	0.428	[0.402,0.452]
$X_3$	$\simeq 0$	[0,0.011]	0.254	[0.232,0.282]
Sobol benchmark				
$X_1$	0.759	[0.701,0.813]	0.769	[0.718,0.807]
$X_2$	0.146	[0.123,0.175]	0.290	[0.258,0.316]
$X_3$	0.025	[0.017,0.034]	0.038	[0.022,0.053]
$X_4$	0.003	[0,0.010]	0.019	[0.010,0.029]
$X_5$ à $X_8$	$\simeq 0$	$\simeq 0$	$\simeq 0$	$\simeq 0$

TABLE 1 – Indice de sensibilité de premier ordre et totaux pour le modèle linéaire, d'I higami et de Sobol.

Le indice de sensibilité de premier ordre et totaux, estimés à l'aide de la fonction `sobol2002` du package `sensitivity`, ont donné dans la table 1. La taille d'échantillon utilisée est 10000, et les intervalles de confiance ont été obtenus par 100 répliques bootstrap.

**Lecture des résultats** Nous présentons ci-dessous un exemple de lecture des résultats pour le modèle d'I higami :

- La variable qui a le plus d'influence sur la variance de la sortie (au sens de l'indice total, c'est-à-dire en prenant en compte les interactions avec les autres variables), est la variable  $X_1$ , avec un indice total de 0.6 et près de 30% de la variance de  $Y$  expliquée à elle seule.
- La variable  $X_2$  n'intervient que seule (indice d'ordre un équivalent à indice total), en expliquant près de 40% de la variance de  $Y$ .
- La variable  $X_3$  n'a aucune influence seule, mais a une influence relativement importante en interaction (avec  $X_1$ ), avec un indice total d'environ 0.3.
- On en déduit que la variance de  $Y$  est due pour 40% à  $X_2$ , 30% à  $X_1$  et 30% à l'interaction entre  $X_1$  et  $X_3$ .

Notons également que dans cet exemple, l'interaction entre  $X_1$  et  $X_3$  est due à une relation non additive entre ces deux variables dans l'expression du modèle.

**Interprétation des résultats** Afin d'interpréter la valeur de l'indice de sensibilité, nous fixons tour à tour chaque variable du modèle à son espérance et examinons l'impact que cela a sur la distribution de la sortie  $Y$ . La figure 2 présente sous la forme de boîte à moustache le résultat obtenu, pour les trois modèles linéaire, d'Ishigami et de Sobol (de gauche à droite). Sur chaque graphique, la première boîte (à gauche) correspond à la distribution initiale de  $Y$ . Les boîtes suivantes correspondent aux distributions de  $Y$  lorsque la variable est fixée une à une, en les ordonnant de gauche à droite selon leur ordre décroissant d'importance (au sens de l'indice de sensibilité total).

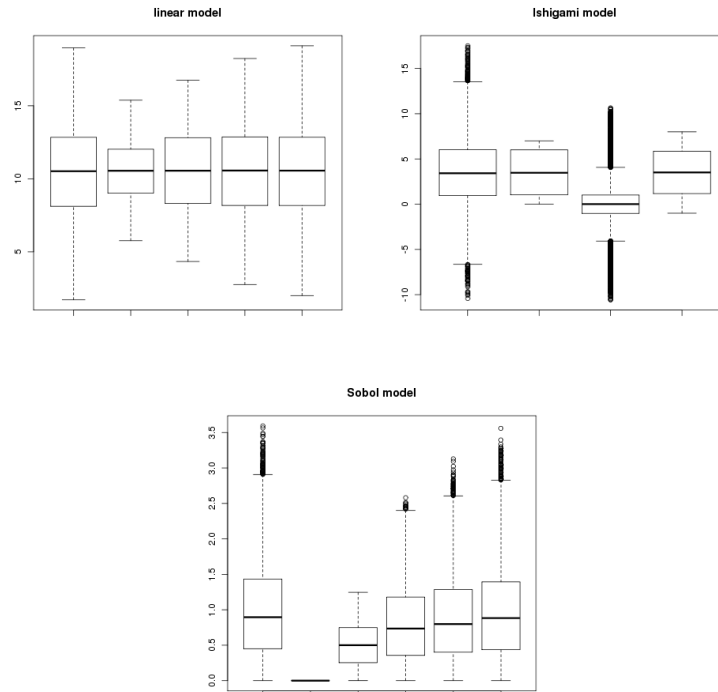


FIGURE 2 – Distribution de  $Y$  en fonction de la variable d'entrée fixée, pour les trois modèles linéaire, d'Ishigami et de Sobol.

Comme attendu, la plus grande réduction de variance est obtenue en fixant la variable ayant l'indice de sensibilité total le plus important. Mais il faut avoir à l'esprit que modifier la distribution d'une variable d'entrée (pour réduire sa variance) n'agit pas uniquement sur la variance de  $Y$  : en effet, au cas d'un modèle linéaire, une modification de la distribution de l'entrée influence également la position centrale de la distribution. Il conviendra donc de s'assurer avant de modifier la distribution d'une variable d'entrée dans le but d'améliorer le pouvoir prédictif du modèle, que celle-ci est bien justifiée et réaliste.



### 4.3 Illustration de modèles à entrées dépendantes

Dans cette seconde illustration numérique, nous considérons le modèle :

$$Y = X_1 + X_2 X_3 + X_4^2;$$

où  $(X_1; X_2; X_3; X_4)^t$  est un vecteur gaussien d'espérance  $(1; 1; 1; 1)^t$  et de matrice de variance

$$\Sigma = \begin{pmatrix} 4 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 4 & 2 \\ 0 & 0 & 2 & 3 \end{pmatrix};$$

L'indice de sensibilité multidimensionnel de premier ordre et total ont été estimés à l'aide de la fonction `sobol_multi` du package `sensitivity-dependent` avec une taille d'échantillon de 10000. Les résultats (estimation moyenne et écart-type sur 100 répétitions bootstrap) sont donnés par la table 2 et la figure 3. La table 2 présente également les résultats d'estimation de l'indice d'ordre un par la méthode de McKay (20 répétitions de l'échantillonnage LHS), obtenus par la fonction `mckay` du package `sensitivity-dependent`.

variable	indice ordre 1 (McKay)	indice multidimensionnel ordre 1	total
$X_1$	0.073	0.061 (0.014)	0.069 (0.008)
$X_2$	0.060	0.049 (0.013)	0.298 (0.016)
$X_3$	0.078	0.634 (0.016)	0.882 (0.012)
$X_4$	0.526		

TABLE 2 – Indices de sensibilité de premier ordre et totaux du modèle d'Ishigami.

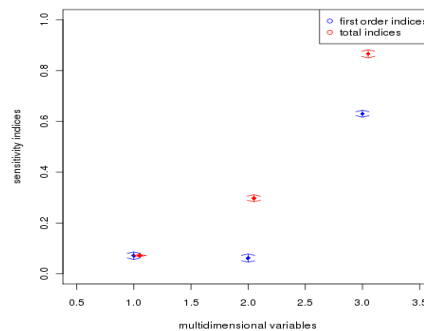


FIGURE 3 – Indices de sensibilité multidimensionnel de premier ordre et totaux

Sur cet exemple, les indices d'ordre un estimés par McKay nous indiquent une importance prépondérante de la variable  $X_4$ , et aucune influence de l'autre

variable (eule). Le calcul de indice multidimensionnel nous permet d'aller plus loin dans l'interprétation :

- la variable  $X_2$  a également une influence significative grâce à de interaction avec d'autre variable ,
- le variable  $X_3$  et  $X_4$  (et non pas uniquement  $X_4$ ) expliquent certes à elle seule une grande partie de la variance de  $Y$  (environ 60%), mais également une autre partie non négligeable de cette variance à travers l'interaction avec la variable  $X_2$ .

## 5 Discussion

L'analyse de sensibilité globale a pour objectif de déterminer l'impact de variable d'entrée sur la variabilité de la sortie d'un modèle mathématique. Dans le cas de modèle à entrée indépendante (cas le plus fréquemment abordé dans la littérature mais pas forcément le plus répandu en pratique), les indices de sensibilité expriment la part de variance de la sortie due à chaque variable d'entrée. De nombreux travaux ont permis de développer de méthode d'estimation efficace et simple à mettre en oeuvre, que le praticien pourra intégrer dans son problème à propre code de calcul. Nous attirons néanmoins l'attention du praticien quant à l'interprétation et l'utilisation de résultats de sensibilité, comme nous l'avons illustré précédemment : modifier la variance de la variable d'entrée la plus influente pour diminuer l'incertitude de prédiction du modèle n'influe pas uniquement sur la variance de la sortie. L'hypothèse d'indépendance de entrée faite précédemment est primordiale et le praticien ne doit surtout pas s'aventurer à de analyse de sensibilité classique lorsque son modèle ne respecte pas cette hypothèse. Dans une telle situation, il doit plutôt se servir de indice de sensibilité multidimensionnel lorsque les entrées ne sont pas toutes dépendantes entre elles, soit de indice de sensibilité d'ordre un classique mais estimé par de méthode spécifique au cas d'entrée dépendante : méthode de McKay, facile d'implémentation, ou de Da Veiga, plus complexe mais beaucoup moins gourmande en évaluation du modèle.

## Références

- [1] S. Da Veiga, F. Wahl, and F. Gamboa. Local polynomial estimation for sensitivity analysis on model with correlated input. *Technometrics*, 51(4) :452–463, 2009.
- [2] T. Homma and A. Saltelli. Importance measure in global sensitivity analysis of non linear model. *Reliability Engineering and System Safety*, 52 :1–17, 1996.
- [3] B. Iou. Revue sur l'analyse de sensibilité globale de modèle numérique. *Journal de la Société Française de Statistique*, 152 :1–23, 2011.

- [4] J. Jacque , C. Lavergne, and N. Devictor. Sensitivity analysis in presence of model uncertainty and correlated input . *Reliability Engineering and System Safety*, 91 :1126–1134, 2006.
- [5] M.D. McKay. Evaluating prediction uncertainty. Technical Report NUREG/CR-6311, US Nuclear Regulatory Commission and Los Alamos National Laboratory, 1995.
- [6] M.D. McKay, R. Beckman, and W. Conover. A comparison of three methods for selecting value of input variable in the analysis of output from a computer code. *Technometrics*, 21(2) :239–245, 1979.
- [7] H. Niederreiter. *Random Number Generation and Quasi-Monte Carlo Methods*. Philadelphia : SIAM, 1992.
- [8] A. Owen. *Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing*, chapter Randomly Permuted (t,m, )-Net and (t, )-Sequence . New York : Springer-Verlag, Niederreiter,H. and Shiue,P.J.-S. (Ed ), 1995.
- [9] A. Owen. Monte carlo extension of quasi-monte carlo. In *1998 Winter Simulation Conference*, Washington (DC, USA), 1998.
- [10] G. Pujol and B. Ioo . Package sensitivity : Sensitivity analysis . Technical report, R software, 2008.
- [11] A. Saltelli, K. Chan, and E.M. Scott, editor . *Sensitivity Analysis*. Wiley, 2000.
- [12] A. Saltelli and E.M. Scott. Guest editorial : The role of sensitivity analysis in the corroboration of model and its link to model structural and parametric uncertainty. *Reliability Engineering and System Safety*, 1997.
- [13] A. Saltelli and S. Tarantola. Sensitivity analysis : a prerequisite in model building? *Foresight and Precaution*, 2000.
- [14] A. Saltelli and S. Tarantola. On the relative importance of input factors in mathematical model : safety assessment for nuclear waste disposal. *Journal of the American Statistical Association*, 97(459) :702–709, 2002.
- [15] A. Saltelli, S. Tarantola, and F. Campolongo. Sensitivity analysis as an ingredient of modeling. *Statistical Science*, 15(4) :377–395, 2000.
- [16] I.M. Sobol. Sensitivity estimate for nonlinear mathematical model . *Mathematical Modelling and Computational Experiments*, 1 :407–414, 1993.
- [17] T. Turanyi. Sensitivity analysis of complex kinetic system, tool and application . *Journal of Mathematical Chemistry*, 5 :203–248, 1990.