

GSA of Model Output With Dependent Input: Part I Sampling Techniques

Thierry A. Mara

Joint Research Centre

11th SAMO Summer School, June 6-10 2022, (Online event)

Outline

Introduction

Case 1: Marginal and conditional cdfs known

- The Rosenblatt transformation

- A special case: the Nataf transformation

Case 2: Marginal and conditional pdfs unknown

- Rejection sampling (MCMC)

Introduction

Let $y = \mathcal{M}(\mathbf{x})$ be the scalar model response of interest. When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{x_1} p_{x_2} \dots p_{x_d}$, we saw that it was always possible to turn the problem into the following: $y = f(\mathbf{u})$ with $\mathbf{u} = (u_1, \dots, u_d) \sim \mathcal{U}(0, 1)^d$ by setting $u_i = F_{x_i}(x_i)$.

Introduction

Let $y = \mathcal{M}(\mathbf{x})$ be the scalar model response of interest. When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{x_1} p_{x_2} \dots p_{x_d}$, we saw that it was always possible to turn the problem into the following: $y = f(\mathbf{u})$ with $\mathbf{u} = (u_1, \dots, u_d) \sim \mathcal{U}(0, 1)^d$ by setting $u_i = F_{x_i}(x_i)$.

Case 1: When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{\mathbf{x}} \neq p_{x_1} p_{x_2} \dots p_{x_d}$ it is sometimes possible to turn the problem into $y = f(\mathbf{u})$. But, the transformation is not unique. There are in principle $d!$ possible transformations.

Introduction

Let $y = \mathcal{M}(\mathbf{x})$ be the scalar model response of interest. When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{x_1} p_{x_2} \dots p_{x_d}$, we saw that it was always possible to turn the problem into the following: $y = f(\mathbf{u})$ with $\mathbf{u} = (u_1, \dots, u_d) \sim \mathcal{U}(0, 1)^d$ by setting $u_i = F_{x_i}(x_i)$.

Case 1: When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{\mathbf{x}} \neq p_{x_1} p_{x_2} \dots p_{x_d}$ it is sometimes possible to turn the problem into $y = f(\mathbf{u})$. But, the transformation is not unique. There are in principle $d!$ possible transformations.

But, by knowing the transformation of \mathbf{u} into \mathbf{x} and vice-versa, it is possible to interpret the sensitivity indices of the u -variables as sensitivity indices of the x -variables.

Introduction

Let $y = \mathcal{M}(\mathbf{x})$ be the scalar model response of interest. When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{x_1} p_{x_2} \dots p_{x_d}$, we saw that it was always possible to turn the problem into the following: $y = f(\mathbf{u})$ with $\mathbf{u} = (u_1, \dots, u_d) \sim \mathcal{U}(0, 1)^d$ by setting $u_i = F_{x_i}(x_i)$.

Case 1: When $\mathbf{x} = (x_1, \dots, x_d) \sim p_{\mathbf{x}} \neq p_{x_1} p_{x_2} \dots p_{x_d}$ it is sometimes possible to turn the problem into $y = f(\mathbf{u})$. But, the transformation is not unique. There are in principle $d!$ possible transformations.

But, by knowing the transformation of \mathbf{u} into \mathbf{x} and vice-versa, it is possible to interpret the sensitivity indices of the u -variables as sensitivity indices of the x -variables.

Case 2: Nevertheless, there are situations for which it is not possible to perform such a transformation.

Case 1: Marginal and conditional cdfs known

The Rosenblatt Transform: Suppose known all cdfs:

$$(F_{x_{i_1}}, F_{x_{i_2}|x_{i_1}}, F_{x_{i_3}|x_{i_1}, x_{i_2}}, \dots, F_{x_{i_d}|x_{\sim i_d}}), \forall i_k \in (1, 2, \dots, d).$$

The Rosenblatt Transform: Suppose known all cdfs:

$(F_{x_{i_1}}, F_{x_{i_2}|x_{i_1}}, F_{x_{i_3}|x_{i_1}, x_{i_2}}, \dots, F_{x_{i_d}|x_{\sim i_d}}), \forall i_k \in (1, 2, \dots, d)$. Then one can use the Rosenblatt transform (1952) to sample \mathbf{x} from \mathbf{u} (or the other way around),

$$\begin{cases} x_{i_1} &= F_{x_{i_1}}^{-1}(u_{i_1}) \\ x_{i_2} &= F_{x_{i_2}}^{-1}(u_{i_2}|u_{i_1}) \\ \vdots &\vdots \\ x_{i_d} &= F_{x_{i_d}|x_{\sim i_d}}^{-1}(u_{i_d}|\mathbf{u}_{\sim i_d}) \end{cases} \quad (1)$$

N.B.: **The Rosenblatt transformation is not unique** as there are $d!$ possible transformations.

Example

We want a sample of $(x_1, x_2) \in \mathcal{U}(0, 1)^2$ uniformly distributed over the triangle $x_1 + x_2 \leq 1$. The problem being symmetric, we have

$$F_{x_1} = F_{x_2} \text{ and } F_{x_1|x_2} = F_{x_2|x_1}.$$

Example

We want a sample of $(x_1, x_2) \in \mathcal{U}(0, 1)^2$ uniformly distributed over the triangle $x_1 + x_2 \leq 1$. The problem being symmetric, we have

$$F_{x_1} = F_{x_2} \text{ and } F_{x_1|x_2} = F_{x_2|x_1}.$$

It can be shown that the marginal cdf is:

$$F_{x_i}(x_i) = 1 - (1 - x_i)^2 = u_i$$

and the conditional cdf is: $F_{x_j|x_i}(x_i, x_j) = \frac{x_j}{1-x_i} = u_j$

Example

We want a sample of $(x_1, x_2) \in \mathcal{U}(0, 1)^2$ uniformly distributed over the triangle $x_1 + x_2 \leq 1$. The problem being symmetric, we have

$$F_{x_1} = F_{x_2} \text{ and } F_{x_1|x_2} = F_{x_2|x_1}.$$

It can be shown that the marginal cdf is:

$$F_{x_i}(x_i) = 1 - (1 - x_i)^2 = u_i$$

and the conditional cdf is: $F_{x_j|x_i}(x_i, x_j) = \frac{x_j}{1-x_i} = u_j$

To sample \mathbf{x} from a sample of \mathbf{u} we must make the following transformations,

$$\begin{cases} x_i &= 1 - \sqrt{1 - u_i} \\ x_j &= u_j \sqrt{1 - u_i} \end{cases} \quad (i, j) = (1, 2) \text{ or } (i, j) = (2, 1)$$

Cholesky transformation

From $\boldsymbol{u} \sim \mathcal{U}(0, 1)^d$
to
 $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

Let $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be a random vector of RVs normally distributed. $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$ is the vector of means and $\boldsymbol{\Sigma}$ is a $d \times d$ Covariance Matrix (symmetric & positive-definite). If $\boldsymbol{\Sigma}$ is diagonal, then the RVs are **independent** otherwise they are **correlated**.

Let $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be a random vector of RVs normally distributed. $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$ is the vector of means and $\boldsymbol{\Sigma}$ is a $d \times d$ Covariance Matrix (symmetric & positive-definite). If $\boldsymbol{\Sigma}$ is diagonal, then the RVs are **independent** otherwise they are **correlated**.

We note that

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T} \quad (2)$$

with $|\cdot|$ is the determinant.

Let $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be a random vector of RVs normally distributed. $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$ is the vector of means and $\boldsymbol{\Sigma}$ is a $d \times d$ Covariance Matrix (symmetric & positive-definite). If $\boldsymbol{\Sigma}$ is diagonal, then the RVs are **independent** otherwise they are **correlated**.

We note that

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T} \quad (2)$$

with $|\cdot|$ is the determinant.

By setting $\mathbf{z} = (\mathbf{x} - \boldsymbol{\mu})\mathbf{U}^{-1}$ where \mathbf{U} is the **upper triangular Cholesky matrix** defined as $\boldsymbol{\Sigma} = \mathbf{U}^T \mathbf{U}$, Eq.(2) becomes,

$$\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}_d) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\mathbf{z}\mathbf{z}^T} \quad (3)$$

which means that \mathbf{z} is a vector of **independent standard normal variables**.

Let $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be a random vector of RVs normally distributed. $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$ is the vector of means and $\boldsymbol{\Sigma}$ is a $d \times d$ Covariance Matrix (symmetric & positive-definite). If $\boldsymbol{\Sigma}$ is diagonal, then the RVs are **independent** otherwise they are **correlated**.

We note that

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T} \quad (2)$$

with $|\cdot|$ is the determinant.

By setting $\mathbf{z} = (\mathbf{x} - \boldsymbol{\mu})\mathbf{U}^{-1}$ where \mathbf{U} is the **upper triangular Cholesky matrix** defined as $\boldsymbol{\Sigma} = \mathbf{U}^T \mathbf{U}$, Eq.(2) becomes,

$$\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}_d) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\mathbf{z}\mathbf{z}^T} \quad (3)$$

which means that \mathbf{z} is a vector of **independent standard normal variables**. We get $\mathbf{x} = \boldsymbol{\mu} + \mathbf{z}\mathbf{U} \Rightarrow$ Samples of \mathbf{x} can be generated from samples of \mathbf{z} which can be generated from \mathbf{u} .

Cholesky transformation

Exercises

Exercise 1: We want a sample of $\mathcal{N}\left(\mathbf{x} \mid \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}\right)$. Set $N = 128$, and generate a sample of $(u_1, u_2) \sim \mathcal{U}(0, 1)^2$.

1. Transform u_1 into $z_1 \sim \mathcal{N}(z_1|0, 1)$
2. Transform u_2 into $z_2 \sim \mathcal{N}(z_2|0, 1)$
3. Find \mathbf{U} the upper Cholesky matrix of the covariance matrix
4. Deduce a sample of \mathbf{x} . Check the empirical covariance matrix.

From $\mathbf{u} \sim \mathcal{U}(0, 1)^d$ to
 $\mathbf{x} \sim c \cdot p_{x_1} \cdot p_{x_2} \cdot p_{x_3} \cdots p_{x_d}$
with c a Gaussian copula density

Nataf transformation

The Nataf Transform: Let $\mathbf{x} = (x_1, \dots, x_d)$ be a random vector of correlated RVs distributed w.r.t. $(F_{x_1}, \dots, F_{x_d})$ with correlation matrix \mathbf{C}_{xx} .

Nataf transformation

The Nataf Transform: Let $\mathbf{x} = (x_1, \dots, x_d)$ be a random vector of correlated RVs distributed w.r.t. $(F_{x_1}, \dots, F_{x_d})$ with correlation matrix \mathbf{C}_{xx} . To generate samples of \mathbf{x} , Iman & Conover (1982) have proposed the following algorithm

1. Let $\mathbf{u} \sim \mathcal{U}(0, 1)^d$
2. Transform \mathbf{u} into $\tilde{\mathbf{x}}$ with the integral transform method, that is, $\tilde{x}_i = F_{x_i}^{-1}(u_i)$
3. Transform \mathbf{u} into $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{C}_{xx})$ with the Cholesky transformation
4. Transform $\tilde{\mathbf{x}}$ into \mathbf{x} such that $\text{rank}(x_i) = \text{rank}(z_i)$,
 $\forall i = 1, \dots, d$

N.B.: Iman & Conover's method ensures that \mathbf{x} and \mathbf{z} has the same Rank (Spearman) Correlation Matrix. If the target is the (Pearson) Correlation Matrix, then one might need to modify \mathbf{C}_{xx} in Step 3. In that case, the technique is known as the Nataf Transform (1962).

About the Rank Transformation:

Let consider the following samples, Their ranks are resp.,

$$\tilde{x}_1 = \begin{bmatrix} 0.55 \\ 0.31 \\ 0.78 \\ 0.03 \\ 0.27 \end{bmatrix}, z_1 = \begin{bmatrix} -0.15 \\ 0.61 \\ 0.38 \\ -0.31 \\ 0.91 \end{bmatrix}$$

$$\text{rank}(\tilde{x}_1) = \begin{bmatrix} 4 \\ 3 \\ 5 \\ 1 \\ 2 \end{bmatrix}, \text{rank}(z_1) = \begin{bmatrix} 2 \\ 4 \\ 3 \\ 1 \\ 5 \end{bmatrix}$$

About the Rank Transformation:

Let consider the following samples, Their ranks are resp.,

$$\tilde{x}_1 = \begin{bmatrix} 0.55 \\ 0.31 \\ 0.78 \\ 0.03 \\ 0.27 \end{bmatrix}, z_1 = \begin{bmatrix} -0.15 \\ 0.61 \\ 0.38 \\ -0.31 \\ 0.91 \end{bmatrix} \quad \text{rank}(\tilde{x}_1) = \begin{bmatrix} 4 \\ 3 \\ 5 \\ 1 \\ 2 \end{bmatrix}, \text{rank}(z_1) = \begin{bmatrix} 2 \\ 4 \\ 3 \\ 1 \\ 5 \end{bmatrix}$$

They would have the same rank, for instance, by rearranging \tilde{x}_1 as

$$\text{follows, } x_1 = \begin{bmatrix} 0.27 \\ 0.55 \\ 0.31 \\ 0.03 \\ 0.78 \end{bmatrix} \rightarrow \text{rank}(x_1) = \begin{bmatrix} 2 \\ 4 \\ 3 \\ 1 \\ 5 \end{bmatrix}$$

Exercises

Exercise: We want a sample of (x_1, x_2) with $p_{x_1} = \mathcal{U}(1, 2)$, $p_{x_2} = \mathcal{N}(0, 2)$ and (Pearson) Correlation Matrix

$\mathbf{C}_{xx} = \begin{bmatrix} 1 & -0.7 \\ -0.7 & 1 \end{bmatrix}$. Set $N = 128$, and generate a sample of $(u_1, u_2) \sim \mathcal{U}(0, 1)^2$. Set $\mathbf{R}_{xx} = \mathbf{C}_{xx}$

1. Transform the sample of \mathbf{u} into a sample of $\tilde{\mathbf{x}} \sim p_{x_1} \cdot p_{x_2}$
2. Transform the sample of \mathbf{u} into a sample of $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{R}_{xx})$
3. Rank-transform $\tilde{\mathbf{x}}$ to obtain \mathbf{x}
4. Check the (Pearson) correlation matrix of \mathbf{x} . If not satisfactory change \mathbf{R}_{xx} and go to Step 2.

With Rosenblatt we can perform the following transformation:

$$\mathbf{u} \rightarrow \mathbf{x}$$

We Nataf transformation we do: $\mathbf{u} \rightarrow \mathbf{z} \rightarrow \mathbf{x}$

In the end of the day we can write:

$$y = \mathcal{M}(\mathbf{x}) = f(\mathbf{u})(= g(\mathbf{z}))$$

With Rosenblatt we can perform the following transformation:

$$\mathbf{u} \rightarrow \mathbf{x}$$

We Nataf transformation we do: $\mathbf{u} \rightarrow \mathbf{z} \rightarrow \mathbf{x}$

In the end of the day we can write:

$$y = \mathcal{M}(\mathbf{x}) = f(\mathbf{u})(= g(\mathbf{z}))$$

With the pick-freeze method (seen on Monday morning) to compute the Sobol' indices how would you proceed:

$$S_i^{IA} = \frac{2 \sum_{n=1}^N \left(y_n^A - y_n^{A_{b_i}} \right) \left(y_n^{B_{a_i}} - y_n^B \right)}{\sum_{n=1}^N \left(y_n^A - y_n^B \right)^2 + \left(y_n^{A_{b_i}} - y_n^{B_{a_i}} \right)^2}$$
$$T_i^{IA} = \frac{\sum_{n=1}^N \left(y_n^A - y_n^{A_{b_i}} \right)^2 + \left(y_n^{B_{a_i}} - y_n^B \right)^2}{\sum_{n=1}^N \left(y_n^A - y_n^B \right)^2 + \left(y_n^{A_{b_i}} - y_n^{B_{a_i}} \right)^2}$$

With Rosenblatt we can perform the following transformation:

$$\mathbf{u} \rightarrow \mathbf{x}$$

We Nataf transformation we do: $\mathbf{u} \rightarrow \mathbf{z} \rightarrow \mathbf{x}$

In the end of the day we can write:

$$y = \mathcal{M}(\mathbf{x}) = f(\mathbf{u})(= g(\mathbf{z}))$$

How would you generate the output samples: $(\mathbf{y}^A, \mathbf{y}^B, \mathbf{y}^{A_{b_i}}, \mathbf{y}^{B_{a_i}})$?

With Rosenblatt we can perform the following transformation:

$$\mathbf{u} \rightarrow \mathbf{x}$$

We Nataf transformation we do: $\mathbf{u} \rightarrow \mathbf{z} \rightarrow \mathbf{x}$

In the end of the day we can write:

$$y = \mathcal{M}(\mathbf{x}) = f(\mathbf{u})(= g(\mathbf{z}))$$

Answer: Generate $(\mathbf{u}^A, \mathbf{u}^B, \mathbf{u}^{A_{b_i}}, \mathbf{u}^{B_{a_i}})$, transform them into $(\mathbf{x}^A, \mathbf{x}^B, \mathbf{x}^{A_{b_i}}, \mathbf{x}^{B_{a_i}})$, run the model and collect the model responses which correspond to $(\mathbf{y}^A, \mathbf{y}^B, \mathbf{y}^{A_{b_i}}, \mathbf{y}^{B_{a_i}})$.

With Rosenblatt we can perform the following transformation:

$$\mathbf{u} \rightarrow \mathbf{x}$$

We Nataf transformation we do: $\mathbf{u} \rightarrow \mathbf{z} \rightarrow \mathbf{x}$

In the end of the day we can write:

$$y = \mathcal{M}(\mathbf{x}) = f(\mathbf{u})(= g(\mathbf{z}))$$

But how do we interpret the sensitivity indices of the u -variables?

Besides, RT is not unique, which RT shall we use?

Answer in Part II.

Cholesky transformation

Compute the variance-based sensitivity indices of the linear function:

$$y = \sum_{i=1}^3 x_i$$

with $p(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$

where $\boldsymbol{\mu} = (1, 2, 3)$ and

$$\mathbf{C} = \begin{array}{ccc} & \begin{matrix} x_1 & x_2 & x_3 \end{matrix} \\ \begin{bmatrix} 1 & -0.5 & 0 \\ -0.5 & 1 & 0.8 \\ 0 & 0.8 & 1 \end{bmatrix} & \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \end{array} \quad \mathbf{C} = \begin{array}{ccc} & \begin{matrix} x_2 & x_3 & x_1 \end{matrix} \\ \begin{bmatrix} 1 & 0.8 & -0.5 \\ 0.8 & 1 & 0 \\ -0.5 & 0 & 1 \end{bmatrix} & \begin{matrix} x_2 \\ x_3 \\ x_1 \end{matrix} \end{array}$$

Generate the samples of \mathbf{u} with $N = 128$. Compute the Sobol' indices (S_i^{IA}, T_i^{IA}) with the RT (here Cholesky transformation) of (u_1, u_2, u_3) into (x_1, x_2, x_3) . What-if we RT transform (u_2, u_3, u_1) into (x_2, x_3, x_1) ?

Case 2: Marginal and conditional cdfs unknown

Rejection Sampling: Let \mathbf{x} be a random vector of RVs distributed w.r.t. the joint pdf $p_{\mathbf{x}}$. If none of the techniques above can be applied: Try acceptance/rejection sampling techniques like Markov Chains Monte Carlo (MCMC).

But, ignoring the independent u-variables, the sensitivity analysis that can be performed is limited.

Some References

- ▶ M.D. McKay, R.J. Beckman, and W.J. Conover (1979). Technometrics, 239–245.
- ▶ Nataf, B. (1962), Compt. Rend. Acad. Sci., 42–43
- ▶ Iman, R.I., Conover, W.J. (1982). Commun. Stat. Simul. Comput. 311–334.
- ▶ Rosenblatt, M (1952). Annals Math. Stat., 470–472