

Kernel-based Sensitivity Analysis of set-valued models.
Application to pollutant concentration maps SA
Thesis: Sensitivity Analysis for Constrained Robust Optimization

Noé Fellmann
Céline Helbert & Christophette Blanchet (ECL)
Adrien Spagnol & Delphine Sinoquet (IFP Energies Nouvelles)

École Centrale de Lyon & IFP Énergie nouvelles

SIAM UQ24, MS215 Global Sensitivity Analysis and Feature Importance

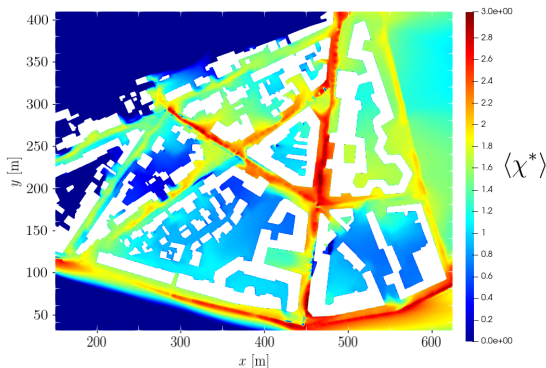


Model of pollutant concentration maps [M. Pasquier]

Input $\mathbf{U} = (\theta, U_\infty, q, \beta, \nu_{\max})$ with,

- Wind direction $\theta \sim \mathcal{N}_{[5,10]}(0.7, 0.5^2)$ [rad]
- Wind speed $U_\infty \sim \mathcal{N}_{[-0.35;1.75]}(8, 2)$ [m/s]
- Traffic volume $q \sim \mathcal{SN}_{[100;500]}(450, 100, -3)$ [vehicle/h]
- Proportion of diesel and petrol engine $\beta \sim \mathcal{U}([0, 1])$ [-]
- Speed limit $\nu_{\max} \sim \mathcal{U}([30; 50])$ [km/h]

Output : $(x,y) \mapsto \Phi_{\mathbf{U}}(x,y)$ a pollutant concentration map



Sensitivity Analysis of map-valued models

Map-valued model :

$$\Phi : \mathbf{U} \mapsto \Phi_{\mathbf{U}}$$

where

$$\Phi_{\mathbf{U}} : (x, y) \mapsto \Phi_{\mathbf{U}}(x, y) \in \mathbb{R}$$

Goal : Quantify the effect of the inputs \mathbf{U} on the spatial output $\Phi_{\mathbf{U}}$

Sensitivity Analysis of map-valued models

Map-valued model :

$$\Phi : \mathbf{U} \mapsto \Phi_{\mathbf{U}}$$

where

$$\Phi_{\mathbf{U}} : (x, y) \mapsto \Phi_{\mathbf{U}}(x, y) \in \mathbb{R}$$

Goal : Quantify the effect of the inputs \mathbf{U} on the spatial output $\Phi_{\mathbf{U}}$

Sensitivity analysis context

$$(U_1, \dots, U_d) \xrightarrow{f} Z = f(U_1, \dots, U_d)$$

How can the uncertainty of Z be divided and allocated to the uncertainty of the inputs U_i ?

- Sobol indices : $S_i = \frac{\text{Var} \mathbb{E}(Z|U_i)}{\text{Var} Z}$
- Dependence measures : $S_i = \|\mathbb{P}_{(U_i, Z)} - \mathbb{P}_{U_i} \otimes \mathbb{P}_Z\|$

Screening : U_1, \dots, U_k are influential and U_{k+1}, \dots, U_d are not influential

Ranking : $U_1 \prec \dots \prec U_d$

Table of Contents

- 1 Pointwise Sensitivity Analysis of pollutant concentration maps
- 2 Set-valued models
- 3 Sensitivity Analysis with kernel-based indices
- 4 Kernel-based Sensitivity Analysis for sets

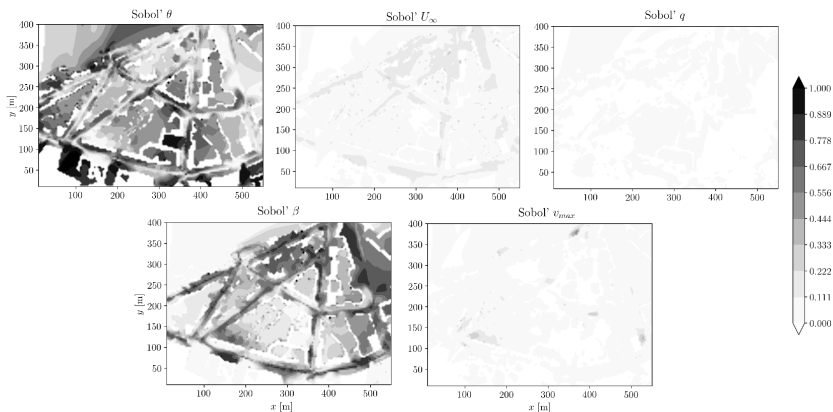
Table of Contents

- 1 Pointwise Sensitivity Analysis of pollutant concentration maps
- 2 Set-valued models
- 3 Sensitivity Analysis with kernel-based indices
- 4 Kernel-based Sensitivity Analysis for sets

Sobol indices maps [M. Pasquier]

Sobol indices at the position (x, y) :

$$S_i(x, y) = \frac{\text{Var}[\mathbb{E}[\Phi_{\mathbf{U}}(x, y) | U_i]]}{\text{Var}[\Phi_{\mathbf{U}}(x, y)]}$$



Aggregated Sobol indices [M. Pasquier]

Aggregated Sobol' indices (Gamboa et al. 2013) :

$$S_i^{\text{gen}} := \sum_{j=1}^m w_j S_i(x^{(j)}, y^{(j)}) \quad \text{with} \quad w_j = \frac{\text{Var}[\Phi_{\mathbf{U}}(x^{(j)}, y^{(j)})]}{\sum_{k=1}^m \text{Var}[\Phi_{\mathbf{U}}(x^{(k)}, y^{(k)})]}.$$

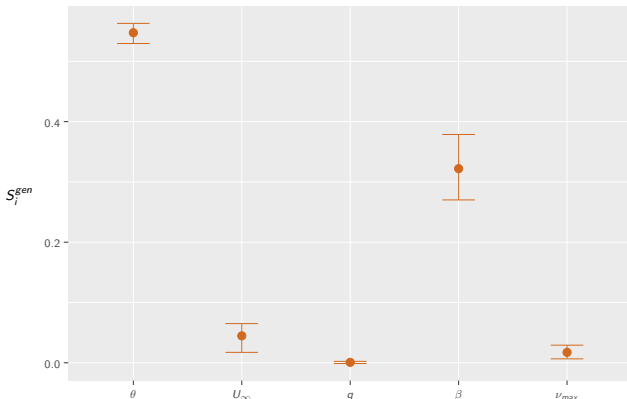


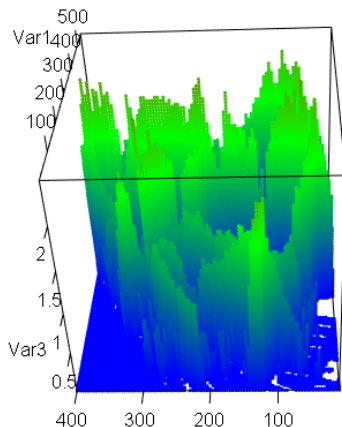
Figure – Estimated aggregated first-order Sobol' indices with 2^{12} model evaluations.

Table of Contents

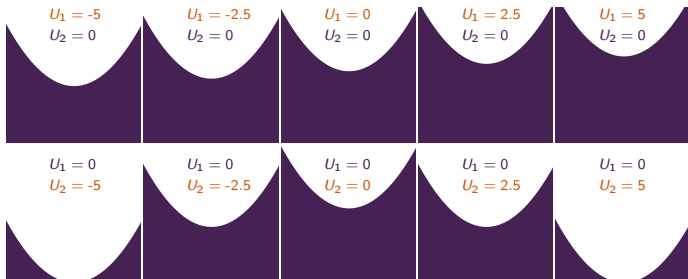
- 1 Pointwise Sensitivity Analysis of pollutant concentration maps
- 2 Set-valued models
- 3 Sensitivity Analysis with kernel-based indices
- 4 Kernel-based Sensitivity Analysis for sets

From map-valued model to set-valued model

$$\Psi : \begin{array}{l} \mathcal{U} \rightarrow \\ \mathbf{u} \mapsto \end{array} \begin{array}{l} \mathcal{L}(\mathcal{X}) \\ \Gamma_{\mathbf{u}} = \{(x, y, c) \in \mathcal{D} \times [0, C_{max}], c \leq \Phi_{\mathbf{u}}(x, y)\} \end{array}$$



Sensitivity analysis of set-valued models?



How to do sensitivity analysis of set-valued models?

Sensitivity analysis on the volume

How to do sensitivity analysis of set-valued models?

- Conduct sensitivity analysis on the volume : $U \rightarrow \lambda(\Gamma_U)$

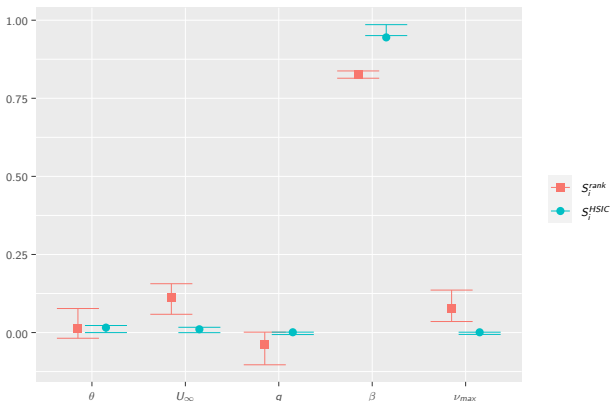


Figure – Estimation of Sobol indices (rank-method) and HSIC-based indices (rbf_anova kernel) of the volume of Γ_U . 1000 model evaluations are used and confidence interval are estimated with 100 bootstrap resamples.

Sensitivity analysis on the volume

How to do sensitivity analysis of set-valued models?

- Conduct sensitivity analysis on the volume : $U \rightarrow \lambda(\Gamma_U)$

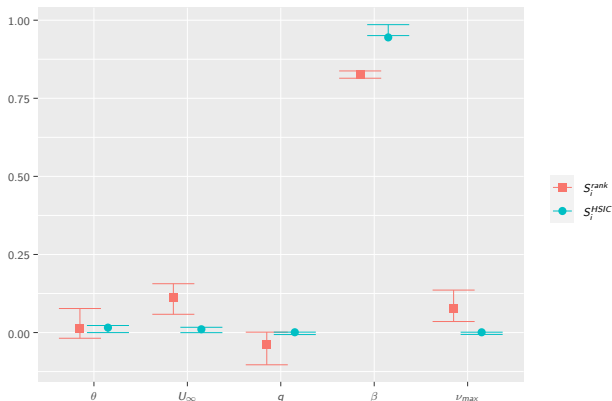


Figure – Estimation of Sobol indices (rank-method) and HSIC-based indices (rbf_anova kernel) of the volume of Γ_U . 1000 model evaluations are used and confidence interval are estimated with 100 bootstrap resamples.

Table of Contents

- 1 Pointwise Sensitivity Analysis of pollutant concentration maps
- 2 Set-valued models
- 3 Sensitivity Analysis with kernel-based indices
- 4 Kernel-based Sensitivity Analysis for sets

Distribution embedding into a RKHS

Given a kernel $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \mapsto \mathbb{R}$,

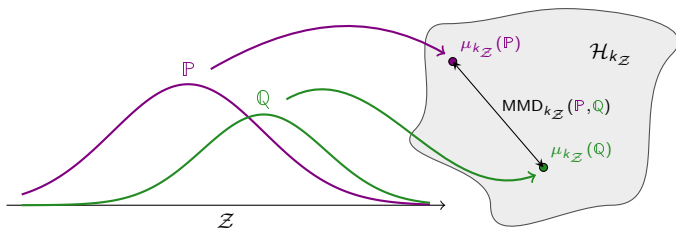


Figure – Kernel mean embedding

- $\mu_{k_{\mathcal{Z}}}(\mathbb{P}) = \int_{\mathcal{Z}} k(\cdot, z) d\mathbb{P}(z)$
- $\text{MMD}_{k_{\mathcal{Z}}}(\mathbb{P}, \mathbb{Q}) = \|\mu_{k_{\mathcal{Z}}}(\mathbb{P}) - \mu_{k_{\mathcal{Z}}}(\mathbb{Q})\|_{\mathcal{H}_{k_{\mathcal{Z}}}}$

The MMD is a distance between distribution iff $k_{\mathcal{Z}}$ is a characteristic kernel, i.e. the mean embedding is injective.

From MMD to HSIC

Dependence measure $S_i = \|\mathbb{P}_{(U_i, Z)} - \mathbb{P}_{U_i} \otimes \mathbb{P}_Z\|$ to quantify the effect of an input U_i on the entire output distribution.

Hilbert Schmidt Independence Criterion (HSIC), Gretton et al. 2006

With $K = k_{U_i} \otimes k_Z$, the HSIC is given by :

$$\begin{aligned} \text{HSIC}_{K_i}(U_i, Z) &:= \text{MMD}_{K_i}^2(\mathbb{P}_{U_i, Z}, \mathbb{P}_{U_i} \mathbb{P}_Z) \\ &= \|\mu_{K_i}(\mathbb{P}_{U_i, Z}) - \mu_{k_{U_i}}(\mathbb{P}_{U_i}) \otimes \mu_{k_Z}(\mathbb{P}_Z)\|_{\mathcal{H}_K}^2 \end{aligned}$$

If k_{U_i} and k_Z are **characteristic**, then

$$\text{HSIC}_{K_i}(U_i, Z) = 0 \text{ iff } U_i \perp Z$$

Screening is then possible with independence testing.

HSIC-ANOVA indices [daVeiga 2021]

Assuming that the inputs are **independent** and that the input kernels are **ANOVA**,

$$\text{HSIC}(\mathbf{U}, Z) = \sum_{A \subseteq \{1, \dots, d\}} \sum_{B \subseteq A} (-1)^{|A|-|B|} \text{HSIC}(\mathbf{U}_B, Z).$$

HSIC-ANOVA indices are then defined as :

$$S_i^{\text{HSIC}} := \frac{\text{HSIC}(U_i, Z)}{\text{HSIC}(\mathbf{U}, Z)},$$

$$S_{T_i}^{\text{HSIC}} := 1 - \frac{\text{HSIC}(\mathbf{U}_{-i}, Z)}{\text{HSIC}(\mathbf{U}, Z)}$$

and are suited for **ranking** (and screening).

- Easy to estimate :

$$\text{HSIC}(\mathbf{U}_A, Z) = \mathbb{E} [(K_A(\mathbf{U}_A, \mathbf{U}_A') - 1)k_Z(Z, Z')].$$

- Only requirement : to have kernels on the inputs and on the output

Table of Contents

- 1 Pointwise Sensitivity Analysis of pollutant concentration maps
- 2 Set-valued models
- 3 Sensitivity Analysis with kernel-based indices
- 4 Kernel-based Sensitivity Analysis for sets

HSIC ANOVA indices for sets, definition of the indices

\mathcal{Z}	\longleftrightarrow	$\mathcal{L}(\mathcal{X})$ the space of Lebesgue measurable subsets of \mathcal{X}
Z	\longleftrightarrow	Γ a random set
$k_{\mathcal{Z}}$	\longleftrightarrow	k_{set} a (characteristic) kernel on $\mathcal{L}(\mathcal{X})$

Proposition (Fellmann, Blanchet-Scalliet et al. 2023)

The function k_{set} defined by :

$$k_{\text{set}}(\gamma_1, \gamma_2) = \exp\left(-\frac{\lambda(\gamma_1 \Delta \gamma_2)}{2\sigma^2}\right)$$

- is a bounded and measurable kernel
- is characteristic.

$$\text{HSIC}_{K_i \otimes k_{\text{set}}}(U_i, \Gamma) = \|\mu_{K_i \otimes k_{\text{set}}}(\mathbb{P}_{(U_i, \Gamma)}) - \mu_{K_i \otimes k_{\text{set}}}(\mathbb{P}_{U_i} \otimes \mathbb{P}_{\Gamma})\|_{\mathcal{H}_{K_i \otimes k_{\text{set}}}}^2$$

HSIC ANOVA indices for sets, estimation

$$\text{HSIC}(\mathbf{U}_A, Z) = \mathbb{E} \left[(K_A(\mathbf{U}_A, \mathbf{U}_A') - 1) k_{\text{set}}(\Gamma, \Gamma') \right].$$

The indices can be estimated using :

$$\widehat{H}_{\text{set}}(U_i, \Gamma) = \frac{2}{n(n-1)} \sum_{j < l}^n \left(K_i(U_i^{(j)}, U_i^{(l)}) - 1 \right) \exp\left(-\frac{\lambda(\mathcal{X})}{2\sigma^2} \hat{\lambda}_m(\Gamma^{(j)} \Delta \Gamma^{(l)})\right).$$

Input kernels :

- the Sobolev kernel of order 1,
 $k_{\text{sob}}(x, y) = 1 + (x - \frac{1}{2})(y - \frac{1}{2}) + \frac{1}{2}[(x - y)^2 - |x - y| + \frac{1}{6}]$
- the Gaussian kernel, $k_{\text{rbf}}(x, y) = e^{-\frac{1}{2}(\frac{x-y}{\sigma})^2}$ with $\sigma > 0$,
- the Laplace kernel, $k_{\text{exp}}(x, y) = e^{-\frac{|x-y|}{h}}$ with $h > 0$,
- the Matérn 3/2, $k_{3/2}(x, y) = \left(1 + \sqrt{3} \frac{|x-y|}{h}\right) e^{-\sqrt{3} \frac{|x-y|}{h}}$ with $h > 0$,
- the Matérn 5/2, $k_{5/2}(x, y) = \left(1 + \sqrt{5} \frac{|x-y|}{h} + \frac{5}{3} \frac{|x-y|}{h^2}\right) e^{-\sqrt{5} \frac{|x-y|}{h}}$ with $h > 0$.

HSIC ANOVA indices for sets, results

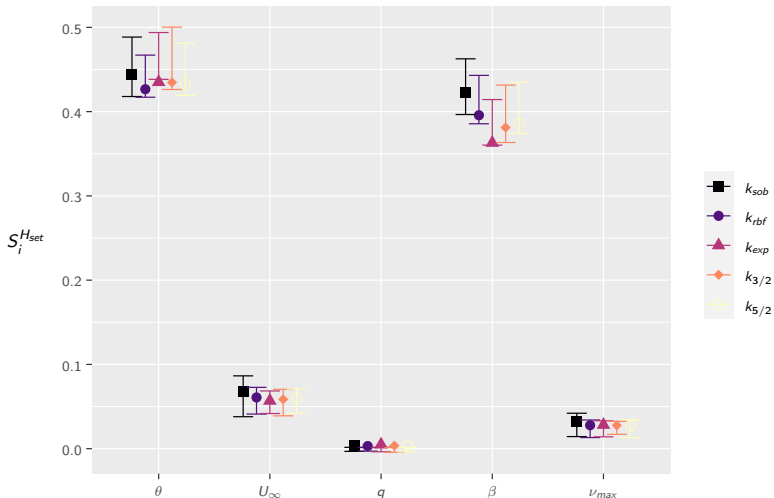


Figure – Estimation of $S_i^{H_{set}}$ for five input kernels, 1000 model evaluations. Confidence intervals are obtained by bootstrap with 100 resamples

Comparison with others indices

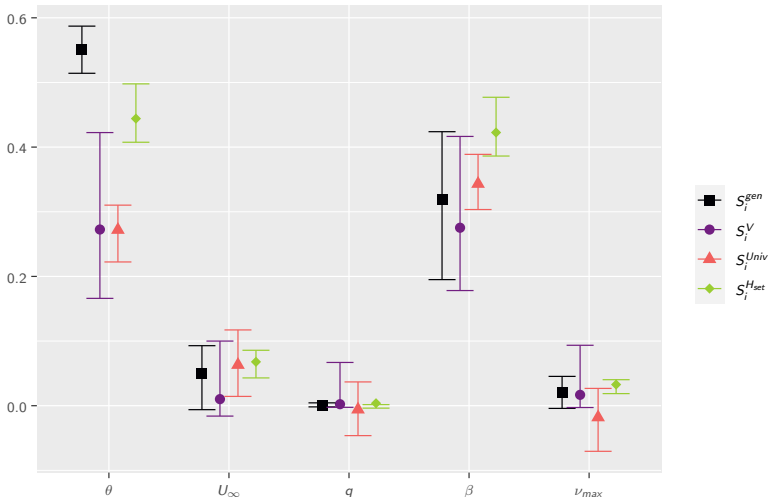


Figure – Comparison of the four indices with a total budget of $n = 1000$ model evaluations. 100 bootstrap sample are used to estimate confidence intervals [Fellmann, Pasquier et al. 2023]

Comparison with other indices

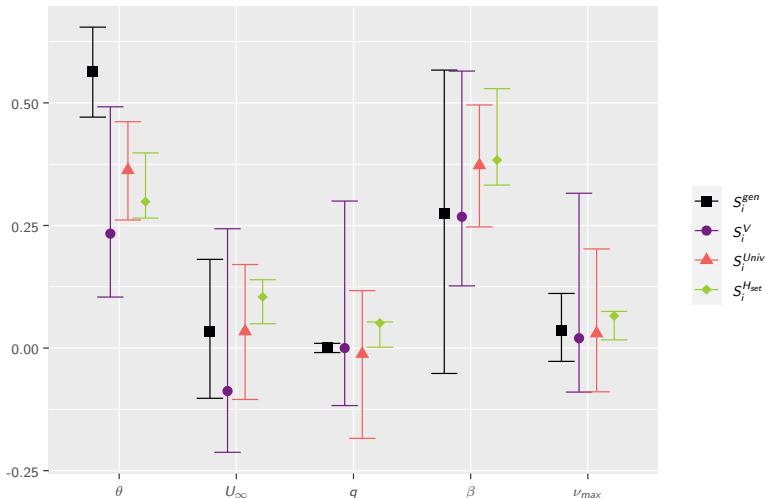


Figure – Comparison of the four indices with a total budget of $n = 100$ model evaluations. 100 bootstrap sample are used to estimate confidence intervals [Fellmann, Pasquier et al. 2023]

Conclusion






- HSIC-ANOVA indices are
 - ✓ cheap to estimate
 - ✓ suited for screening
 - ✓ suited for ranking
 - ✗ sometimes hard to interpret (interactions)
 - ✓ very permissive in the type of the output (and of the inputs?)
- The presented methodology for map-valued outputs but can be used for any set-valued outputs.

Conclusion

- HSIC-ANOVA indices are
 - ✓ cheap to estimate
 - ✓ suited for screening
 - ✓ suited for ranking
 - ✗ sometimes hard to interpret (interactions)
 - ✓ very permissive in the type of the output (and of the inputs?)
- The presented methodology for map-valued outputs but can be used for any set-valued outputs.

Conclusion

- HSIC-ANOVA indices are
 - ✓ cheap to estimate
 - ✓ suited for screening
 - ✓ suited for ranking
 - ✗ sometimes hard to interpret (interactions)
 - ✓ very permissive in the type of the output (and of the inputs?)
- The presented methodology for map-valued outputs but can be used for any set-valued outputs.
- Thanks!

-  daVeiga, Sébastien (jan. 2021). "Kernel-based ANOVA decomposition and Shapley effects - Application to global sensitivity analysis". [working paper or preprint. url : https://hal.archives-ouvertes.fr/hal-03108628](https://hal.archives-ouvertes.fr/hal-03108628).
-  Fellmann, Noé, Christophe Blanchet-Scalliet et al. (2023). [Kernel-based sensitivity analysis for \(excursion\) sets. arXiv : 2305.09268 \[math.ST\]](https://arxiv.org/abs/2305.09268).
-  Fellmann, Noé, Mathis Pasquier et al. (2023). "Sensitivity analysis for sets : application to pollutant concentration maps". In : [arXiv preprint](https://arxiv.org/abs/2305.09268).
-  Gamboa, Fabrice et al. (2013). [Sensitivity analysis for multidimensional and functional outputs. arXiv : 1311.1797 \[stat.AP\]](https://arxiv.org/abs/1311.1797).
-  Gretton, Arthur et al. (2006). "A Kernel Method for the Two-Sample-Problem". In : [Advances in Neural Information Processing Systems. T. 19. MIT Press. url : https://proceedings.neurips.cc/paper/2006/hash/e9fb2eda3d9c55a0d89c98d6c54b5b3e-Abstract.html](https://proceedings.neurips.cc/paper/2006/hash/e9fb2eda3d9c55a0d89c98d6c54b5b3e-Abstract.html).