

DE LA RECHERCHE À L'INDUSTRIE



# Data-driven kernel representations for sampling with an unknown block dependence

**G. Perrin**<sup>(1)</sup>,

in collaboration with C. Soize<sup>(2)</sup> and N.  
Ouhbi<sup>(3)</sup>

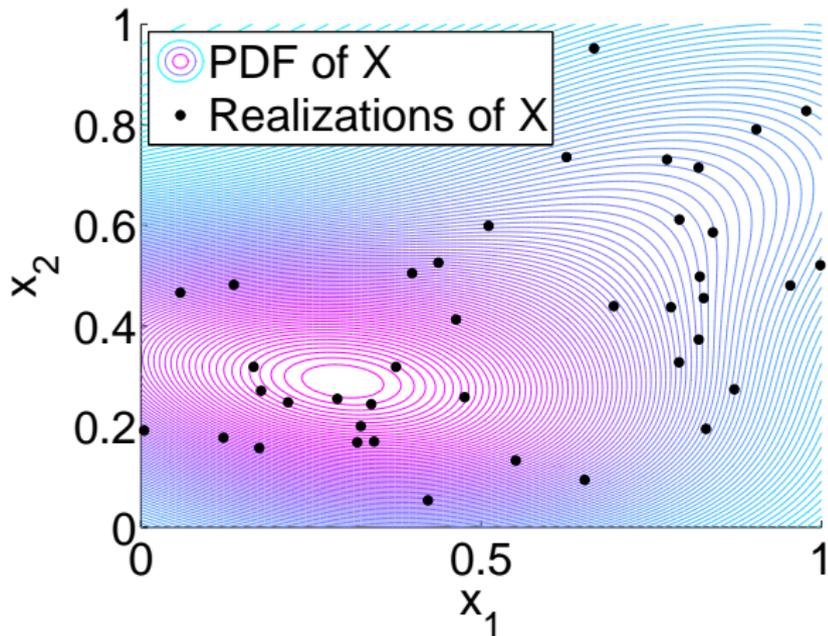
<sup>(1)</sup>CEA/DAM/DIF, Arpajon, France

<sup>(2)</sup>Université Paris-Est, Marne-la-Vallée, France

<sup>(3)</sup>Innovation and Research Department, SNCF, Paris, France

MascotNum 2018 |

March, 22<sup>th</sup> 2018



- When the maximal information about a  $d$ -dimensional random vector  $\mathbf{X}$  is a set of  $N$  iid realizations, the **kernel density estimation** (KDE) is a widely used technique to infer the PDF of  $\mathbf{X}$ .

- For instance, it can be used
  - to process the samples provided by a MCMC approach in the Bayesian calibration of a computational code [Berliner, 2001, Kaipio and Somersalo, 2004],
  - to approximate goal-oriented Sobol indices [Perrin and Defaux, 2018],
  - to optimize under uncertainties a particular code output [Soize and Ghanem, 2017].
- In practice, this technique is limited to cases when  $d$  is **small** (less than five in general).

## Problematic

What could we propose to extend the validity of this technique to higher values of  $d$  ( $d \sim 10 - 100$ ) with limited information ( $N \sim 10d$  for instance)?

- 1 Introduction
- 2 Kernel representations for statistical inference
- 3 Data-driven tensor-product representation in high dimension
- 4 Identification of the mechanical properties of a random medium
- 5 Conclusion

- Let  $\mathbf{X} := \{\mathbf{X}(\omega), \omega \in \Omega\}$  be a second-order random vector defined on a probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ , with values in  $\mathbb{R}^d$ .
- The PDF of  $\mathbf{X}$  is denoted by  $p_{\mathbf{X}}$ .
- The maximal available information about  $p_{\mathbf{X}}$  is a set of  $N > d$  independent and distinct realizations of  $\mathbf{X}$ , which are gathered in the deterministic set  $\mathcal{S}(N) := \{\mathbf{X}(\omega_n), 1 \leq n \leq N\}$ .
- Given these realizations of  $\mathbf{X}$ , the kernel estimator of  $p_{\mathbf{X}}$  is

$$\hat{p}_{\mathbf{X}}(\mathbf{x}; \mathbf{H}, \mathcal{S}(N)) = \frac{\det(\mathbf{H})^{-1/2}}{N} \sum_{n=1}^N K\left(\mathbf{H}^{-1/2}(\mathbf{x} - \mathbf{X}(\omega_n))\right),$$

where  $\det(\cdot)$  is the determinant operator,  $K$  is any function of  $\mathcal{M}_1(\mathbb{R}^d, \mathbb{R}^+)$ , and  $\mathbf{H}$  is a  $(d \times d)$ -dimensional positive definite symmetric "**bandwidth matrix**".

- In this work, we focus on the classical case when  $K$  is the Gaussian multidimensional density ( $\leftrightarrow$  "G-KDE"):

$$\hat{p}_{\mathbf{X}}(\mathbf{x}; \mathbf{H}, \mathcal{S}(N)) = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}; \mathbf{X}(\omega_n), \mathbf{H}), \quad \mathbf{x} \in \mathbb{R}^d.$$

- Here,  $\phi(\cdot; \boldsymbol{\mu}, \mathbf{C})$  is the PDF of any  $\mathbb{R}^d$ -dimensional Gaussian random vector with mean  $\boldsymbol{\mu}$  and covariance matrix  $\mathbf{C}$ :

$$\phi(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C}) := \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{(2\pi)^{d/2} \sqrt{\det(\mathbf{C})}}, \quad \mathbf{x} \in \mathbb{R}^d.$$

$\Rightarrow p_{\mathbf{X}}$  is approximated by a mixture of  $N$  Gaussian PDFs, whose means are the available realizations of  $\mathbf{X}$  and whose covariance matrices are all equal to  $\mathbf{H}$

- By construction,  $\mathbf{H}$  characterizes the **local contribution** of each realization of  $\mathbf{X}$ .
- Its value has to be optimized to minimize the difference between  $p_{\mathbf{X}}$ , which is unknown, and  $\hat{p}_{\mathbf{X}}(\cdot; \mathbf{H}, \mathcal{S}(N))$ .
- The mean integrated squared error (MISE) performance criterion

$$\text{MISE}(\mathbf{H}; d, N) = \mathbb{E} \left[ \int_{\mathbb{R}^d} (p_{\mathbf{X}}(\mathbf{x}) - \hat{p}_{\mathbf{X}}(\mathbf{x}; \mathbf{H}, \mathcal{S}(N)))^2 d\mathbf{x} \right]$$

is generally considered to quantify such a difference.

- Given sufficient regularity conditions on  $p_{\mathbf{X}}$ , an asymptotic approximation of this criterion can be derived, leading to the commonly-used **Silverman bandwidth matrix** [Silverman, 1986]

$$\mathbf{H}^{\text{Silv}}(d, N) := (h^{\text{Silv}}(d, N))^2 \hat{\mathbf{R}}_{\mathbf{X}}, \quad h^{\text{Silv}}(d, N) := \left( \frac{1}{N} \frac{4}{(d+2)} \right)^{\frac{1}{d+4}},$$

with  $\hat{\mathbf{R}}_{\mathbf{X}}$  the empirical estimation of the covariance matrix of  $\mathbf{X}$ .

- In practice, it is generally observed that, for fixed values of  $N$ ,  $\mathbf{H}^{\text{Silv}}(d, N)$  **overestimates** the scattering of  $p_{\mathbf{X}}$ .
- To circumvent this problem, the LOO expression of the likelihood,

$$\mathcal{L}^{\text{LOO}}(\mathcal{S}(N)|\mathbf{H}) := \prod_{n=1}^N \frac{1}{N-1} \sum_{m=1, m \neq n}^N \phi_{n,m}(\mathbf{H}),$$

$$\phi_{n,m}(\mathbf{H}) := \phi(\mathbf{X}(\omega_n); \mathbf{X}(\omega_m), \mathbf{H}), \quad 1 \leq n, m \leq N,$$

can instead directly be used to identify  $\mathbf{H}$  [van der Laan et al., 2004].

- In this presentation, we will focus on the **maximum likelihood estimate** of  $\mathbf{H}$ , denoted by

$$\mathbf{H}^{\text{MLE}}(d, N) := \arg \max_{\mathbf{H} \in \mathbb{M}^+(d)} \mathcal{L}^{\text{LOO}}(\mathcal{S}(N)|\mathbf{H}).$$

- Considering that the best available approximations of the true mean and covariance matrix of  $\mathbf{X}$  are given by their empirical estimations, the former expression can be slightly modified.
- Indeed, if the PDF of  $\mathbf{X}$  is equal to

$$\tilde{p}_{\mathbf{X}}(\cdot; \mathbf{H}, \mathcal{S}(N)) := \frac{1}{N} \sum_{n=1}^N \phi(\cdot; \mathbf{A}\mathbf{X}(\omega_n) + \boldsymbol{\beta}, \mathbf{H}),$$

$$\boldsymbol{\beta} := (\mathbf{I}_d - \mathbf{A})\hat{\boldsymbol{\mu}}, \quad \mathbf{H} := \hat{\mathbf{R}}_{\mathbf{X}} - \frac{N-1}{N} \mathbf{A}\hat{\mathbf{R}}_{\mathbf{X}}\mathbf{A}^T,$$

the mean and the covariance matrix of  $\mathbf{X}$  are equal to  $\hat{\boldsymbol{\mu}}$  and  $\hat{\mathbf{R}}_{\mathbf{X}}$  respectively [Perrin et al., 2018].

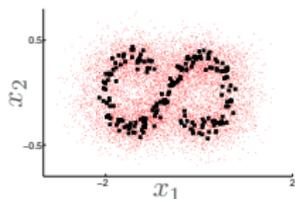
- Given  $\mathcal{S}(N)$ , the G-KDE of the PDF of  $\mathbf{X}$  under constraints on its mean and its covariance matrix will be denoted by  $\tilde{p}_{\mathbf{X}}(\cdot; \mathbf{H}^{\text{MLE}}(d, N), \mathcal{S}(N))$  in the following.

- In practice, when considering the nonparametric modelling of high dimensional random vectors ( $d \sim 10 - 100$ ) with limited information ( $N \sim 10d$  for instance), we observe that  $\mathbf{H}^{\text{MLE}}(d, N)$  is very close to  $\widehat{\mathbf{R}}_X$ .
- This means that we are approximating the PDF of  $\mathbf{X}$  as a **unique** Gaussian PDF, whose parameters correspond to the empirical mean and covariance matrix of  $\mathbf{X}$ :

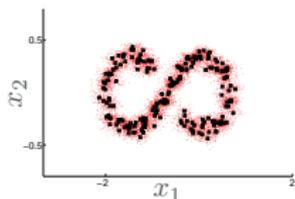
$$\lim_{\mathbf{H} \rightarrow \widehat{\mathbf{R}}_X} \tilde{p}_X(\cdot; \mathbf{H}, \mathcal{S}(N)) = \phi(\cdot; \widehat{\boldsymbol{\mu}}, \widehat{\mathbf{R}}_X).$$

- This could prevent us from recovering the subset of  $\mathbb{R}^d$  on which  $\mathbf{X}$  is actually concentrated.

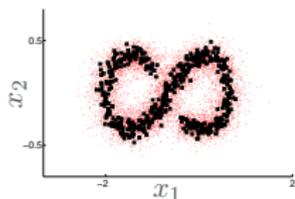
$\mathbf{X}^{(2D)} = (X_1^L + 0.05\xi_1, X_2^L + 0.05\xi_2, \xi_3, \dots, \xi_d)$ ,  
 $X^L \leftrightarrow$  Lemniscate function,  $\xi_i$  are iid standard Gaussian r.v.



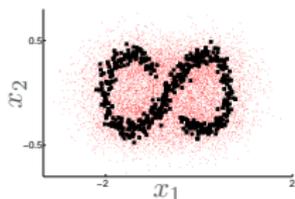
(a)  $d = 2, H^{\text{Silv}}$



(b)  $d = 2, H^{\text{MLE}}$



(c)  $d = 3, H^{\text{MLE}}$

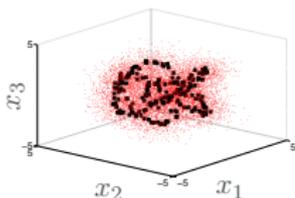


(d)  $d = 5, H^{\text{MLE}}$

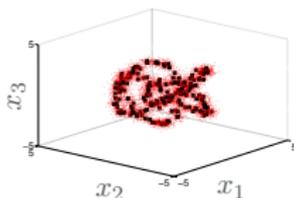
Figure:  $N = 500$  given values of  $\mathbf{X}^{(2D)}$  (big black squares) and  $10^4$  additional values (small red points) generated from a G-KDE approach for several values of  $d$ .

$$\mathbf{X}^{(3D)} = (X_1^{\text{FB}} + \xi_1, \dots, X_3^{\text{FB}} + \xi_3, \xi_4, \dots, \xi_d),$$

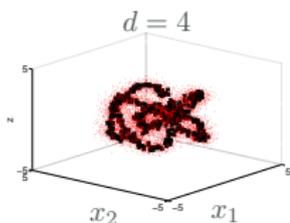
$\mathbf{X}^{\text{FB}} \leftrightarrow$  Four branch clover-knot function,  $\xi_i$  iid standard Gaussian r.v.



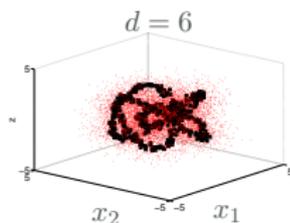
(a)  $d = 3, H^{\text{Silv}}$



(b)  $d = 3, H^{\text{MLE}}$



(c)  $d = 4, H^{\text{MLE}}$



(d)  $d = 6, H^{\text{MLE}}$

Figure:  $N = 500$  given values of  $\mathbf{X}^{(3D)}$  (big black squares) and  $10^4$  additional values (small red points) generated from a G-KDE approach for several values of  $d$ .

- 1 Introduction
- 2 Kernel representations for statistical inference
- 3 Data-driven tensor-product representation in high dimension
- 4 Identification of the mechanical properties of a random medium
- 5 Conclusion

- The idea is to identify groups of components of  $\mathbf{X}$  that can **reasonably** be considered as statistically **independent**, if they exist.
- Given a decomposition of  $\mathbf{X}$  in  $N_b$  blocks  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N_b)}$ , PDF  $p_{\mathbf{X}}$  is approximated as the product of the nonparametric estimations of the PDFs associated with each sub-vector of  $\mathbf{X}$ :

$$p_{\mathbf{X}} \approx \prod_{\ell=1}^{N_b} \tilde{p}_{\mathbf{X}^{(\ell)}}(\cdot; \mathbf{H}^{(\ell)}, \mathcal{S}(N)).$$

- Instead of using statistical tests, we propose to search these groups by looking for the minimum of a cross-validation AIC that is associated with each block formation [Perrin et al., 2018].

## Hypotheses

- 1  $\mathbf{X}$  is **centred** and **uncorrelated**:  $\hat{\boldsymbol{\mu}}_{\mathbf{X}} = \mathbf{0}$ ,  $\hat{\mathbf{R}}_{\mathbf{X}} = \mathbf{I}_d$ .
- 2  $\mathbf{H}_{\ell}$  is parametrized by a **unique scalar**:  $\mathbf{H}_{\ell} = h_{\ell}^2 \mathbf{I}_{d_{\ell}}$ ,  $0 < h_{\ell} \leq 1$ .

- For any  $\mathbf{b}$  in  $\{1, \dots, d\}^d$ ,  $b_i$  can be used as a block index for the  $i^{\text{th}}$  component of  $\mathbf{X}$ .
- This means that if  $b_i = b_j$ ,  $X_i$  and  $X_j$  are supposed to be **dependent** and belong to the **same** block. On the contrary, if  $b_i \neq b_j$ ,  $X_i$  and  $X_j$  are supposed to be **independent** and they can belong to two **different** blocks.
- There exists a bijection between the set of all block by block decompositions of  $\mathbf{X}$  and the set

$$\mathbb{B}(d) := \left\{ \begin{array}{l} \mathbf{b} \in \{1, \dots, d\}^d \mid b_1 = 1, \\ 1 \leq b_j \leq 1 + \max_{1 \leq i \leq j-1} b_i, \quad 2 \leq j \leq d \end{array} \right\}.$$

## Difficulty

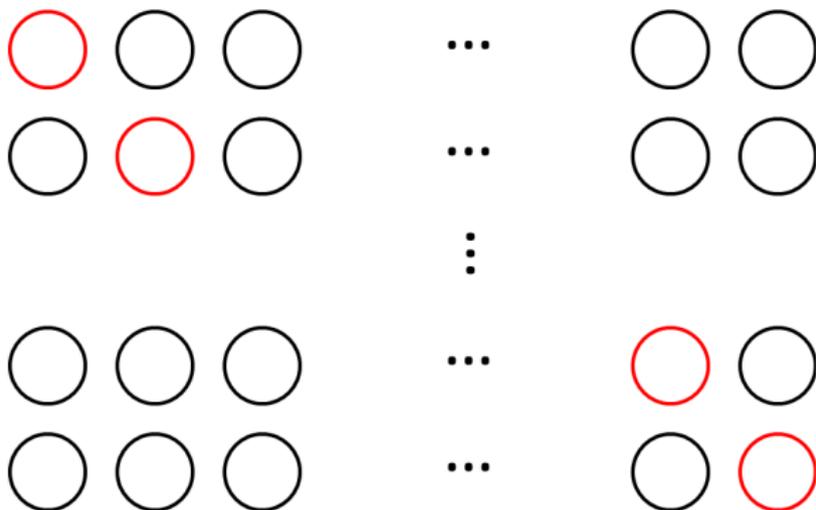
The cardinal of  $\mathbb{B}(d)$  increases exponentially with  $d$ .

## Greedy identification - initialization



$$\mathbf{b} = (1, \dots, 1)$$

## Greedy identification - first loop



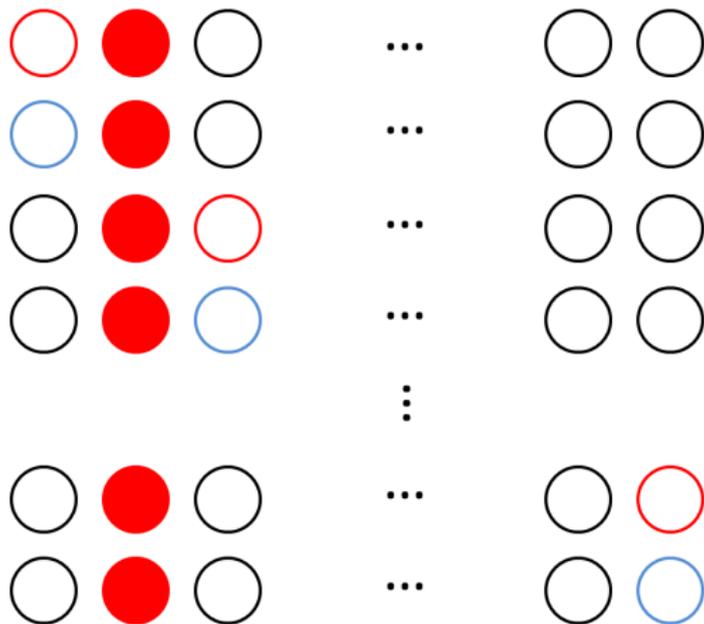
$$\mathbf{b} \in \{(b_1, \dots, b_d) \mid b_i = 2, b_{j \neq i} = 1, 1 \leq i \leq d\}$$

## Greedy identification - first fixed point



$$b = (1, 2, 1, \dots, 1)$$

## Greedy identification - second loop



$$\mathbf{b} \in \{(b_1, \dots, b_d) \mid b_2 = 2, b_i \in \{2, 3\}, b_{j \neq i, j \neq 2} = 1, i \neq 2\}$$

## Greedy identification - second fixed point



$$b = (1, 2, 1, \dots, 1, 3, 1)$$

## Greedy identification - last fixed point

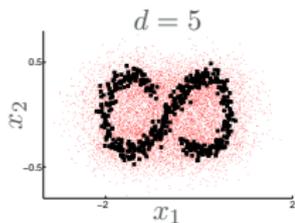


$$b = (1, 2, 3, \dots, 1, 3)$$

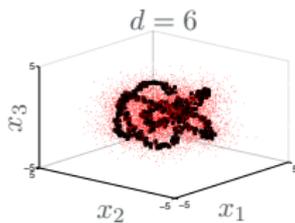
- The greedy algorithm allows the identification of good values for  $\mathbf{b}$  at a reduced computational cost.
- We verified on a series of test cases that the algorithm was able to identify the correct block decomposition with  $2 \leq d \leq 100$  and  $N = 10d$  in a reasonable computational time.
- Evolutionary algorithms can also be used to address problems in higher dimensions.

$d$	1	2	3	4	5	6	7	8	9	10
$\#\mathbb{B}(d)$	1	2	5	15	52	203	877	4140	21147	115975
$N_{\text{greedy}}^{\max}(d)$	1	3	8	17	31	51	78	113	157	211

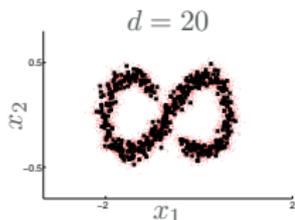
Table: Evolution of the number of elements of  $\mathbb{B}(d)$ ,  $\#\mathbb{B}(d)$ , and the maximum number of configurations tested by the greedy algorithm,  $N_{\text{greedy}}^{\max}(d)$ , with respect to  $d$ .



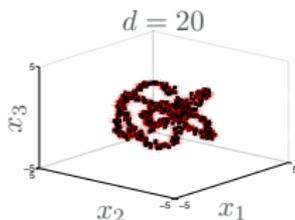
(a)  $\mathbf{b} = (1, \dots, 1)$



(b)  $\mathbf{b} = (1, \dots, 1)$



(c)  $\mathbf{b}$   
 $(1, 1, 2, 3, \dots, 18)$



(d)  $\mathbf{b}$  =  
 $(1, 1, 1, 2, 3, \dots, 17)$

Figure:  $N = 500$  given values of  $\mathbf{X}^{(3D)}$  (big black squares) and  $10^4$  additional values (small red points) generated from a G-KDE approach for several values of  $d$ .

- 1 Introduction
- 2 Kernel representations for statistical inference
- 3 Data-driven tensor-product representation in high dimension
- 4 Identification of the mechanical properties of a random medium
- 5 Conclusion

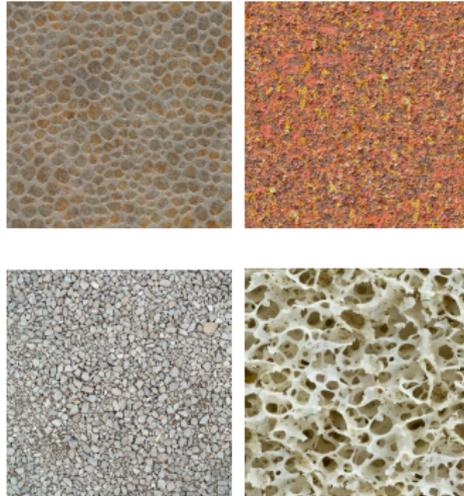


Figure: Four heterogeneous media

## Problematic

How to infer the statistical properties of a random medium from a limited number of indirect measurements?

- We are interested in the identification of the mechanical properties of a **heterogeneous elastic medium**.
- Several experimental tests are performed on a series of  $M$  specimens made of the same material. Let  $\mathcal{V}$  be their volume.
- For each experiment, the applied force field is supposed to be **imposed**, and the induced displacement field is **measured on the contours** of the specimens only.
- In parallel, for given properties of the considered medium, it is possible to approximate (using the Finite Element Method) the displacements that are induced by the experimental force field.
- In this work, we focus on the estimation of 5 quantities gathered in the vector  $\mathbf{z} = (\lambda, \ell_1, \ell_2, \mu_\nu, \mu_E)$ , where  $\lambda$  is a fluctuation level,  $\ell_1$  and  $\ell_2$  are two correlation lengths, and  $\mu_\nu$  and  $\mu_E$  are the means of the Poisson coefficient and the Young modulus respectively.

- Let  $\mathbf{X}$  be the elasticity field characterizing the mechanical properties of the material that constitutes the specimens.
- $\mathbf{X}$  is supposed to be **random**, and we assume that it belongs to a known class of parametric random fields, such that:

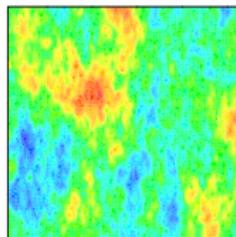
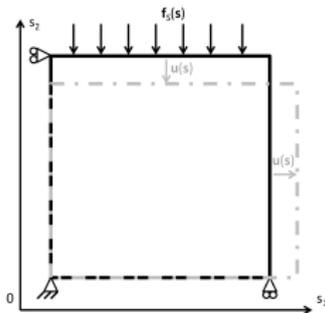
$$\mathbf{X} = \{ \mathbf{X}(\mathbf{s}, \omega; \mathbf{z}^*), \mathbf{s} \in \mathcal{V}, \omega \in \Omega \},$$

where  $\mathbf{z}^* \in \mathbb{Z}$  is unknown.

- $\mathbf{X}$  is not a real-valued random field, but a tensor-valued random field, and its different components cannot be identified separately due to algebraic constraints.
- Let  $\mathbf{u}(\mathbf{X}(\omega; \mathbf{z}))$  be the induced displacement on the contour of the specimen associated with the particular realization  $\mathbf{X}(\omega; \mathbf{z})$  of  $\mathbf{X}(\mathbf{z})$ .
- This displacement can be decomposed in two contributions :

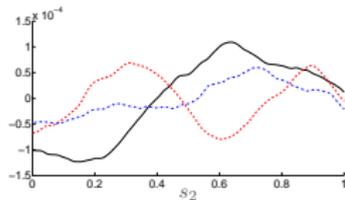
$$\mathbf{u}(\mathbf{s}; \mathbf{X}(\omega; \mathbf{z})) = \hat{\mathbf{u}}(\mathbf{s}; \mathbb{E}[\mathbf{X}(\mathbf{z})]) + \tilde{\mathbf{u}}(\mathbf{s}; \mathbf{X}(\omega; \mathbf{z})).$$

# Illustration of the problem

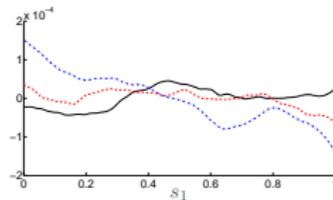


(a) Studied mechanical phenomenon

(b) One possible evolution of the Young modulus



(c) Three measured evolutions of  $\tilde{u}_1(z^*)$



(d) Three measured evolutions of  $\tilde{u}_2(z^*)$

- $z^*$  is modelled by a **random vector**, denoted by  $Z$ , to take into account the fact that its values are unknown. Let  $f_Z$  be its PDF.
- Let  $Y(z)$  be a  $d_y$ -dimensional random vector that condenses the statistical properties of  $u(X(z))$  and  $f_{Y(z)}$  be its PDF.
- $M$  independent realizations of  $Y(z^*)$  are gathered in the set  $\mathbb{Y} := \{Y(\omega_1; z^*), \dots, Y(\omega_M; z^*)\}$  (one for each specimen).
- Using the Bayes theorem, it comes:

$$f_{Z|\mathbb{Y}}(z) = \frac{\mathcal{L}_{\mathbb{Y}}(z)f_Z(z)}{\mathbb{E}[\mathcal{L}_{\mathbb{Y}}(Z)]}, \quad z \in \mathbb{R}^{d_z}.$$

There,  $\mathcal{L}_{\mathbb{Y}}(z)$  is the **likelihood function**.

$$\mathcal{L}_{\mathbf{Y}}(z) = \prod_{m=1}^M f_{\mathbf{Y}(z)}(\mathbf{Y}(\omega_m; z^*)), \quad z \in \mathbb{R}^{d_z}.$$

## Independent estimation

$$\mathcal{L}_{\mathbf{Y}}(z) \approx \prod_{m=1}^M \tilde{f}_{\mathbf{Y}(z)}(\mathbf{Y}(\omega_m; z^*)),$$

where  $\tilde{f}_{\mathbf{Y}(z)}$  is the kernel estimator of  $f_{\mathbf{Y}(z)}$  based on  $N$  independent realizations of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z} = z$ .

For each value of  $z$ ,  $N$  evaluations of the code are required to approximate  $\mathcal{L}_{\mathbf{Y}}(z) \Rightarrow$  as the likelihood function has to be evaluated a high number of times to get precise information about the PDF of  $\mathbf{Z}|\mathbf{Y}$ , this computational cost is generally not affordable.

## Joint estimation

$$\mathcal{L}_Y(z) \approx \prod_{m=1}^M \tilde{f}_{Y(z)}(\mathbf{Y}(\omega_m; \mathbf{z}^*)),$$

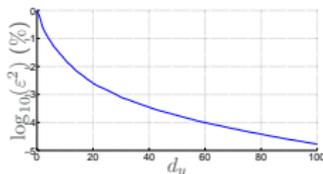
$$\tilde{f}_{Y(z)}(\mathbf{y}) = \frac{\tilde{f}_{Y,Z}(\mathbf{y}, z)}{\int_{\mathbb{R}^{d_y}} \tilde{f}_{Y,Z}(\mathbf{v}, z) d\mathbf{v}}$$

where  $\tilde{f}_{Y,Z}$  is the kernel estimator of the PDF of  $(\mathbf{Y}(\mathbf{Z}), \mathbf{Z})$  based on  $N$  independent realizations of  $(\mathbf{Y}(\mathbf{Z}), \mathbf{Z})$ .

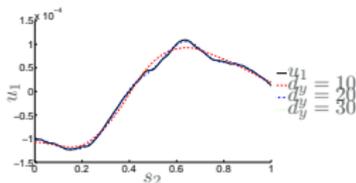
$\Rightarrow$  only  $N$  evaluations of the code are needed to approximate the whole function  $\mathcal{L}_Y$ .

Remark: using Gaussian kernels, the expression of  $\tilde{f}_{Y(z)}$  is **explicit** once  $\tilde{f}_{Y,Z}$  is known.

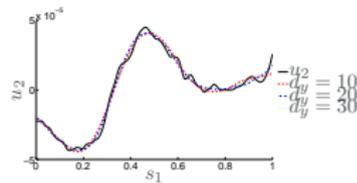
- $N = 1000$  code evaluations are carried out to infer the value of  $z^*$ .
- Coefficients  $\mu_\nu$  and  $\mu_E$  are estimated using preliminary comparisons to the homogeneous case.
- The components of  $\mathbf{Y}(z)$  correspond to the  $d_y$  first components of the KL decomposition of  $\mathbf{u}(\mathbf{X}(z))$ .



(e) Error convergence



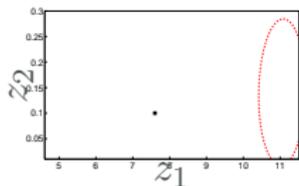
(f) Reduced basis representation of  $u_1(s_2)$



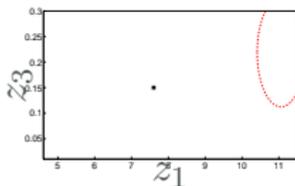
(g) Reduced basis representation of  $u_2(s_1)$

Figure: Evolution of the projection error with respect to  $d_y$ .

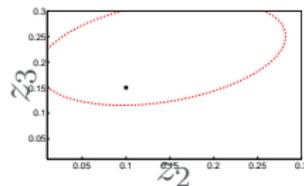
- $N = 1000$  code evaluations are carried out to infer the value of  $z^*$ .
- Coefficients  $\mu_\nu$  and  $\mu_E$  are estimated using preliminary comparisons to the homogeneous case.
- The components of  $\mathbf{Y}(z)$  correspond to the  $d_y$  first components of the KL decomposition of  $\mathbf{u}(\mathbf{X}(z))$ .



(a)  $(-\log(\lambda), \ell_1)$



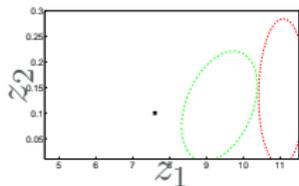
(b)  $(-\log(\lambda), \ell_2)$



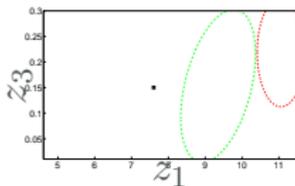
(c)  $(\ell_1, \ell_2)$

Figure: White square: reference value. Red: 95% credible ellipses using  $d_y = 23$  components (corresponding to a projection error of 0.1%).

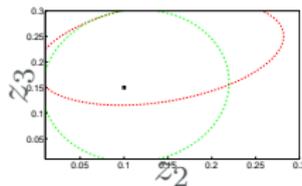
- $N = 1000$  code evaluations are carried out to infer the value of  $z^*$ .
- Coefficients  $\mu_\nu$  and  $\mu_E$  are estimated using preliminary comparisons to the homogeneous case.
- The components of  $\mathbf{Y}(z)$  correspond to the  $d_y$  first components of the KL decomposition of  $\mathbf{u}(\mathbf{X}(z))$ .



(a)  $(-\log(\lambda), \ell_1)$



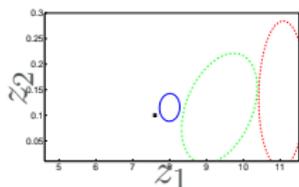
(b)  $(-\log(\lambda), \ell_2)$



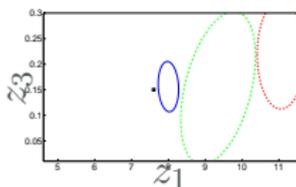
(c)  $(\ell_1, \ell_2)$

Figure: White square: reference value. Red: 95% credible ellipses using  $d_y = 23$  components (corresponding to a projection error of 0.1%). Green: 95% credible ellipses using only  $d_y = 5$  components.

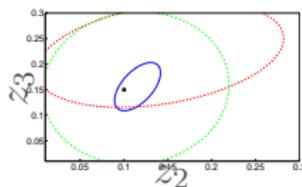
- $N = 1000$  code evaluations are carried out to infer the value of  $\mathbf{z}^*$ .
- Coefficients  $\mu_\nu$  and  $\mu_E$  are estimated using preliminary comparisons to the homogeneous case.
- The components of  $\mathbf{Y}(\mathbf{z})$  correspond to the  $d_y$  first components of the KL decomposition of  $\mathbf{u}(\mathbf{X}(\mathbf{z}))$ .



(a)  $(-\log(\lambda), \ell_1)$



(b)  $(-\log(\lambda), \ell_2)$

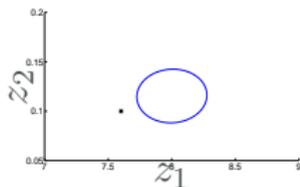


(c)  $(\ell_1, \ell_2)$

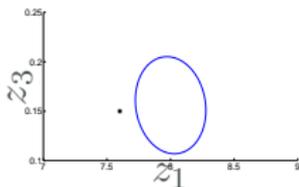
Figure: White square: reference value. Red: 95% credible ellipses using  $d_y = 23$  components (corresponding to a projection error of 0.1%). Green: 95% credible ellipses using only  $d_y = 5$  components. Blue: 95% credible ellipses using  $d_y = 23$  components, with an optimization of the block decomposition of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z}$ .



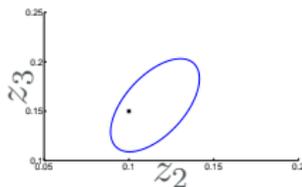
- Step 1:  $N = 1000$  code evaluations are carried out.



(a)  $(-\log(\lambda), \ell_1)$



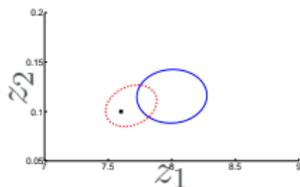
(b)  $(-\log(\lambda), \ell_2)$



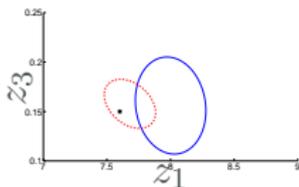
(c)  $(\ell_1, \ell_2)$

Figure: White square: reference value. Blue: results of step 1, with  $d_y = 23$  and 23 blocks for the PDF of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z}$ .

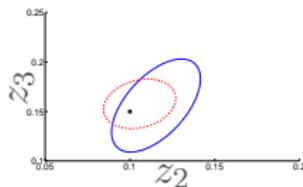
- Step 1:  $N = 1000$  code evaluations are carried out.
- Step 2: 889 new points in the likely region (provided by the calibration results of step 1) of  $z^*$  are added to the learning set.



(a)  $(-\log(\lambda), \ell_1)$



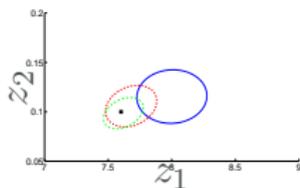
(b)  $(-\log(\lambda), \ell_2)$



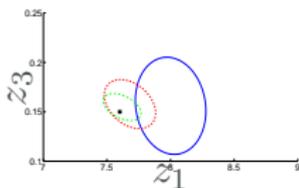
(c)  $(\ell_1, \ell_2)$

Figure: White square: reference value. Blue: results of step 1, with  $d_y = 23$  and 23 blocks for the PDF of  $Y(\mathbf{Z})|\mathbf{Z}$ . Red: results of step 2, with  $d_y = 23$  and 8 blocks for the PDF of  $Y(\mathbf{Z})|\mathbf{Z}$ .

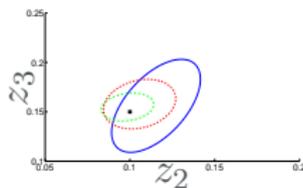
- Step 1:  $N = 1000$  code evaluations are carried out.
- Step 2: 889 new points in the likely region (provided by the calibration results of step 1) of  $z^*$  are added to the learning set.
- Step 3: 631 new points in the likely region (provided by the calibration results of step 2) of  $z^*$  are added to the learning set.



(a)  $(-\log(\lambda), \ell_1)$



(b)  $(-\log(\lambda), \ell_2)$



(c)  $(\ell_1, \ell_2)$

Figure: White square: reference value. Blue: results of step 1, with  $d_y = 23$  and 23 blocks for the PDF of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z}$ . Red: results of step 2, with  $d_y = 23$  and 8 blocks for the PDF of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z}$ . Green: results of step 3, with  $d_y = 23$  and 4 blocks for the PDF of  $\mathbf{Y}(\mathbf{Z})|\mathbf{Z}$ .

- 1 Introduction
- 2 Kernel representations for statistical inference
- 3 Data-driven tensor-product representation in high dimension
- 4 Identification of the mechanical properties of a random medium
- 5 Conclusion

- This work considers the challenging problem of identifying complex PDFs when the maximal available information is a set of independent realizations.
- In that prospect, the multidimensional G-KDE method plays a key role, as it presents a good compromise between complexity and efficiency.
- Two adaptations of this method have been presented to deal with high dimensional random vector:
  - a modified formalism is presented to make the mean and the covariance matrix of the estimated PDF equal to their empirical estimations.
  - tensorized representations are proposed, which are based on the identification of a block by block dependence structure of the random vectors of interest.
- The interest of these two adaptations has been illustrated for the identification of the mechanical properties of a random medium.

- This work considers the challenging problem of identifying complex PDFs when the maximal available information is a set of independent realizations.
- In that prospect, the multidimensional G-KDE method plays a key role, as it presents a good compromise between complexity and efficiency.
- Two adaptations of this method have been presented to deal with high dimensional random vector:
  - a modified formalism is presented to make the mean and the covariance matrix of the estimated PDF equal to their empirical estimations.
  - tensorized representations are proposed, which are based on the identification of a block by block dependence structure of the random vectors of interest.
- The interest of these two adaptations has been illustrated for the identification of the mechanical properties of a random medium.

Thank you for your attention! Questions?

- 

Berliner, L. M. (2001).  
Monte carlo based ensemble forecasting.  
*Statistics and Computing*, 11.
- 

Kaipio, J. P. . and Somersalo, E. (2004).  
*Statistics and Computational Inverse Problems*.  
Springer, New York.
- 

Perrin, G. and Defaux, G. (2018).  
Efficient evaluation of reliability-oriented sensitivity indices.  
*Submitted to Reliability Engineering and System Safety*.
- 

Perrin, G., Soize, C., and Ouhbi, N. (2018).  
Data-driven kernel representations for sampling with an unknown block  
dependence structure under correlation constraints.  
*Journal of Computational Statistics and Data Analysis*, 119:139–154.

-  Silverman, B. W. (1986).  
*Density Estimation for Statistics and Data Analysis*, volume 37.
-  Soize, C. and Ghanem, R. (2017).  
Probabilistic learning on manifold for optimization under uncertainties.  
*Proceeding of Uncecomp 2017*, pages 1–15.
-  van der Laan, M., Dudoit, S., and Keles, S. (2004).  
Asymptotic optimality of likelihood-based cross-validation.  
*Statistical Applications in Genetics and Molecular Biology*, 3(1).

