

DE LA RECHERCHE À L'INDUSTRIE



Nested polynomial trends for the improvement of Gaussian process-based predictors

G. Perrin⁽¹⁾, C. Soize⁽²⁾, J. Garnier⁽³⁾, S. Marque-Pucheu^(1,3)

⁽¹⁾ CEA/DAM/DIF, F-91297, Arpajon, France,

⁽²⁾ Université Paris-Est, MSME UMR 8208 CNRS, Marne-la-Vallée, France

⁽³⁾ Laboratoire de Probabilités et Modèles Aléatoires, Laboratoire Jacques-Louis Lions, Université Paris Diderot, 75205 Paris Cedex 13, France

Ecole thématique ETICS - Barcelonnette |

June 2016

- 1 Introduction
- 2 Gaussian process-based regression (GPR)
- 3 Optimization of the GPR
- 4 Nested polynomial trends
- 5 Conclusions and prospects

- For the last decades, the use of simulation has kept increasing for the analysis of complex and non-linear physical systems.
- Computational models are introduced:
 - to **optimize** the system or its maintenance w.r.t. constraints (geometry, cost, certification criteria...)
 - to **explore** the design space (prototyping, better understanding of the phenomena...)
 - to **evaluate** its robustness and its reliability.



When interested in a complex phenomenon $x \mapsto y(x)$, two kinds of problems can be distinguished:

Forward problems

Given information about x , we would like to:

- 1 **infer** the distribution of x ,
- 2 **compute** some statistical quantities of y , such as:
 - its mean and its variance,
 - probabilities of exceeding threshold,
 - its full density.

Backward problems

Given information about x and y , we would like to:

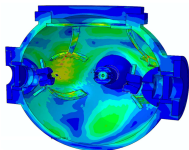
- 1 **Construct** a parametric model $\mathcal{M}(\cdot; \beta)$ for $x \mapsto y(x)$,
- 2 **calibrate** parameters β and validate model \mathcal{M} ,
- 3 **classify** the influence of each input on the variability of y .

The information about x and y is generally a set of experiments (or code evaluations).

Problematic

Most of the methods associated with these problems are based on the post-processing of large sets of simulations. When the numerical cost associated with one simulation is very high (for instance $\sim 100h$ on 100 cores in parallel), we need to know how to

- build **cheap mathematical approximations** of the quantities of interest,
- **control the relevance** of these approximations,
- **use these approximations** to solve the problem we are confronted to (by optimized sequential designs for instance).



- 1 Introduction
- 2 Gaussian process-based regression (GPR)
- 3 Optimization of the GPR
- 4 Nested polynomial trends
- 5 Conclusions and prospects

- Let \mathcal{S} be a physical system, which response depends on a d -dimensional input vector $\mathbf{x} = (x_1, \dots, x_d)$, and which performance can be evaluated from the computation of a quantity of interest, $y(\mathbf{x})$.
- Function y is a **deterministic** mapping that is assumed to be an element of $L^2(\mathcal{D}^d, \mathbb{R})$, where \mathcal{D}^d is a compact set.
- We suppose that the maximal available information about y is a set of N code evaluations, $\{(\mathbf{x}^{(1)}, y(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, y(\mathbf{x}^{(N)}))\}$.
- Each computation of $y(\mathbf{x})$ is computationally **expensive**, such that N is relatively small compared to the complexity of y .
- Given \mathcal{F}_N , the σ -algebra associated with this information, we are interested in the identification of the *best* predictor y^* of y :

$$\|y - y^*\|_{L^2}^2 \leq \|y - \hat{y}\|_{L^2}^2, \quad \hat{y} \in L^2(\mathcal{D}^d, \mathbb{R}).$$

$$(u, v)_{L^2} := \int_{\mathcal{D}^d} u(\mathbf{x})v(\mathbf{x})d\mathbf{x}, \quad \|u\|_{L^2}^2 := (u, u)_{L^2}, \quad u, v \in L^2(\mathcal{D}^d, \mathbb{R}).$$

Hypothesis [Sacks et al., 1989],[Santner et al., 2003]

$x \mapsto y(x)$ is supposed to be a realization of a Gaussian process (GP),

$$Y \sim \text{GP}(\mu, C).$$

If the mean and the covariance functions μ and C are known, then :

$$Y \mid \mathcal{F}_N \sim \text{GP}(\mu_N, C_N),$$

- $\mu_N(x) = \mu(x) + \mathbf{r}(x)^T [C]^{-1} (\mathbb{Y} - \boldsymbol{\mu}),$
- $C_N(x, x') = C(x, x') - \mathbf{r}(x)^T [C]^{-1} \mathbf{r}(x'),$

and $\mathbb{E}[Y \mid \mathcal{F}_N] = \mu_N$ is the **best predictor** of y in the L^2 sense.

$$\mathbb{Y} = (y(x^{(1)}), \dots, y(x^{(N)})), \quad \boldsymbol{\mu} = (\mu(x^{(1)}), \dots, \mu(x^{(N)})),$$

$$[C]_{n,m} = C(x^{(n)}, x^{(m)}), \quad \mathbf{r}(x) = (C(x, x^{(1)}), \dots, C(x, x^{(N)})).$$

Parametric formulation

Assuming that μ and C belong to classes that are parameterized by β and Θ , such that:

$$Y \mid \beta, \Theta \sim \text{GP}(\mu(\beta), C(\Theta)),$$

for all x in \mathcal{D}^d , the best predictor of $y(x)$ is given by:

$$\mathbb{E}[Y(x) \mid \mathcal{F}_N] = \int_{\beta, \Theta, y} y \times \pi[Y(x) = y \mid \beta, \Theta, \mathcal{F}_N] \pi[\beta, \Theta \mid \mathcal{F}_N] d\beta d\Theta dy,$$

where:

- $\pi[Y(x) = y \mid \beta, \Theta, \mathcal{F}_N]$ is Gaussian,
- $\pi[\beta, \Theta \mid \mathcal{F}_N]$ is the **posterior** joint distribution of (β, Θ) .

Linearized plug-in approach [Bect et al., 2012], [Bichon et al., 2008]

- As a good compromise between complexity, efficiency, and errors control, the "linearized plug-in" approach (or "Universal Kriging") is generally preferred to this "full-Bayesian" approach.
- This method consists in :
 - 1 computing the maximum likelihood estimates of β and Θ , which are denoted by β^* and Θ^* ,
 - 2 linearizing the mean function around β^* , such that:

$$\mu(x; \beta) \approx c_0(x) + \langle f(x), \beta \rangle ,$$

$$c_0(x) := \mu(x; \beta^*) - \langle f(x), \beta^* \rangle , \quad f(x) := \frac{\partial \mu}{\partial \beta}(x; \beta^*) ,$$

- 3 reparameterizing: $f(x) \leftarrow (c_0(x), f(x))$, $\beta \leftarrow (1, \beta)$,
- 4 assuming that $\beta \sim \mathcal{U}_{\mathbb{R}^{|\beta|_0}}$, and conditioning all the results by Θ^* .

Linearized plug-in approach [Bect et al., 2012], [Bichon et al., 2008]

It comes:

$$Y \mid \Theta^*, \mathcal{F}_N \sim \text{GP}(\hat{\mu}_N, \hat{C}_N),$$

$$\hat{\mu}_N(\mathbf{x}) = \langle \mathbf{f}(\mathbf{x}), \hat{\boldsymbol{\beta}} \rangle + \mathbf{r}(\mathbf{x})^T [C]^{-1} (\mathbb{Y} - [F] \hat{\boldsymbol{\beta}}),$$

$$\hat{C}_N(\mathbf{x}, \mathbf{x}') = C(\mathbf{x}, \mathbf{x}') - \mathbf{r}(\mathbf{x})^T [C]^{-1} \mathbf{r}(\mathbf{x}') + \mathbf{u}(\mathbf{x})^T ([F]^T [C]^{-1} [F])^{-1} \mathbf{u}(\mathbf{x}').$$

$$[F] = [\mathbf{f}(\mathbf{x}^{(1)}) \cdots \mathbf{f}(\mathbf{x}^{(N)})]^T, \quad \hat{\boldsymbol{\beta}} = ([F]^T [C]^{-1} [F])^{-1} [F]^T [C]^{-1} \mathbb{Y},$$

$$\mathbf{u}(\mathbf{x}) = [F]^T [C]^{-1} \mathbf{r}(\mathbf{x}) - \mathbf{f}(\mathbf{x}).$$

An analytic expression is found back for the best predictor of $y(\mathbf{x})$ in any non-computed point $\mathbf{x} \in \mathcal{D}^d$:

$$\mathbb{E}[Y(\mathbf{x}) \mid \mathcal{F}_N] \approx y^{\text{lin-plug}}(\mathbf{x}) := \hat{\mu}_N(\mathbf{x}).$$

One of the main advantage of this approximated approach comes from the fact that the relevance of the "linearized plug-in" predictor of y can "easily" be computed **without any additional evaluation** of y from **cross-validation** procedures

[Miller, 1974, Dubrulle, 1983, Blatman and Sudret, 2011, Bachoc, 2013]:

$$\left\| y - y^{\text{lin-plug}} \right\|_{L_2}^2 \approx \epsilon_{\text{LOO}}^2 := \frac{1}{N} \sum_{n=1}^N \left(y(\mathbf{x}^{(n)}) - y_{(-n)}^{\text{lin-plug}}(\mathbf{x}^{(n)}) \right)^2,$$

$$y(\mathbf{x}^{(n)}) - y_{(-n)}^{\text{lin-plug}}(\mathbf{x}^{(n)}) = \frac{([R]\mathbb{Y})_n}{[R]_{nn}},$$

$$[R] = [C]^{-1} \{ [I] - [F]([F]^T[C]^{-1}[F])^{-1}[F]^T[C]^{-1} \}.$$

- 1 Introduction
- 2 Gaussian process-based regression (GPR)
- 3 Optimization of the GPR
- 4 Nested polynomial trends
- 5 Conclusions and prospects

To optimize the results associated with such a plug-in approach, two directions can be explored:

- work on the **parameterization of the covariance function** (stationary or not, regularity...),
- work on the **parameterization of the mean function** (linear or not with respect to β , sparse representations...).

- By definition, for all \mathbf{x}, \mathbf{x}' in \mathcal{D}^d :

$$C(\mathbf{x}, \mathbf{x}') := \mathbb{E} \left[(Y(\mathbf{x}) - \mathbb{E}[Y(\mathbf{x})]) \times (Y(\mathbf{x}') - \mathbb{E}[Y(\mathbf{x}')]) \right],$$

such that C is *a priori* any symmetric and non-negative definite kernel that is defined on $\mathcal{D}^d \times \mathcal{D}^d$.

- Remembering that the actual goal is to predict function y , which is defined on \mathcal{D}^d only, it is clear that only "simple" parametric expressions of C can be considered (from a limited number of evaluations of y , we do not pretend to be able to precisely identify C).

A very commonly used representation for C is the following tensorized Matern expression:

$$C(\mathbf{x}, \mathbf{x}') = \sigma^2 \prod_{i=1}^d \frac{1}{2^{\nu-1} \Gamma(\nu)} (2\sqrt{\nu} h_i)^\nu \mathcal{B}_\nu^{III} (2\sqrt{\nu} h_i),$$

where:

- $h_i = g_i \left(\frac{|x_i - x'_i|}{\ell_i} \right)$,
- $\nu \leftrightarrow$ regularity of y ,
- $\sigma^2 \leftrightarrow$ *a priori* uncertainty,
- $\ell_i \leftrightarrow$ correlation lengths,
- $g_i \leftrightarrow$ scaling function (to take into account the fact that the model can be non-stationary with respect to x_i).

- All the parameters associated with covariance function C have to be identified from the available evaluations of y , or be *a priori* fixed (from expert judgment for instance...).
- With only very few information about y and its regularity is available, it is generally accepted that the **Matern-5/2** class is a good *a priori* choice for C , which corresponds to the case when the realizations of Y are twice mean-square differentiable:

$$C(\mathbf{x}, \mathbf{x}') = \sigma^2(1 + \sqrt{h} + (5/3)h^2) \times \exp(-\sqrt{5} h),$$

$$h = \sum_{i=1}^d |x_i - x'_i|/\ell_i.$$

- Fixed parameters: $\nu = 5/2$, $\gamma = \sqrt{5}$, $g(h) = h$,
- Free parameters: σ^2 , $\{\ell_1, \dots, \ell_d\}$.

- By definition, for all \mathbf{x} in \mathcal{D}^d :

$$\mu(\mathbf{x}; \boldsymbol{\beta}) := \mathbb{E} [Y(\mathbf{x})] ,$$

such that μ is *a priori* any function that is defined on \mathcal{D}^d .

- Any parametric expression can therefore be used to characterize μ .
- For instance, one can think to linear or non-linear functions such as polynomials (which can be orthogonal or not in $L^2(\mathcal{D}^d, \mathbb{R})$), cos and sin, neural networks, Low Rank representations...

- Once again, without information about y , polynomials are generally chosen for f . Indeed, the set $\{m_\alpha, \alpha \in \mathbb{N}^d\}$, with

$$m_\alpha(\mathbf{x}) := x_1^{\alpha_1} \times \cdots \times x_d^{\alpha_d}, \quad \mathbf{x} \in L^2(\mathcal{D}^d, \mathbb{R}),$$

defines a basis of $L^2(\mathcal{D}^d, \mathbb{R})$.

- $\beta \mapsto \mu(\beta)$ is therefore chosen **linear**, and for any given value of M , the number of polynomials we want to consider for the representation of μ , the idea is to identify the best M -dimensional subset of $\{m_\alpha, \alpha \in \mathbb{N}^d\}$ to minimize $\|y - \mathbb{E}[Y \mid \mathcal{F}_N]\|_{L^2}$.
- At last, convergence analyses based on the formerly introduced LOO error can be carried out to identify the most accurate value of M .

- In practice, this very complex optimization problem is replaced by an optimization over a **finite** dimensional subset of $\{m_{\alpha}, \alpha \in \mathbb{N}^d\}$.
- Different truncation schemes have thus been proposed to choose such a relevant subset, which are mostly based on the assumption that the most influential elements of $\{m_{\alpha}, \alpha \in \mathbb{N}^d\}$ correspond to the elements of **lowest total polynomial order**:

$$\mathcal{P}(r, d) := \left\{ m_{\alpha} \mid \alpha \in \mathbb{N}^d, \sum_{i=1}^d |\alpha_i| \leq r \right\}.$$

- A **penalization technique**, such as the $\ell - 1$ or the $\ell - 2$ penalizations or the Least Angle Regression (LAR) method [Hastie et al., 2002, Efron et al., 2004, Blatman and Sudret, 2011] are generally used to find these M most significant terms in $\mathcal{P}(r, d)$.
- Cross validation procedures are once again used to avoid overfitting.
- Such an approach will be referred as "LAR+UK" approach in the following [Kersaudy et al., 2015].

Finally, given \mathcal{F}_N , the "LAR+UK" approach follows the following steps:

- 1 identify the M most interesting polynomials to represent y , which are gathered in the vector-valued function $\mathbf{f} = (f_1, \dots, f_M)$, such that:

$$\mu(\boldsymbol{\beta}) := \langle \mathbf{f}, \boldsymbol{\beta} \rangle,$$

- 2 choose an *a priori* adapted parametric expression for $C(\boldsymbol{\Theta})$,
- 3 assume that y is a particular realization of $Y \sim \text{GP}(\mu(\boldsymbol{\beta}), C(\boldsymbol{\Theta}))$,
- 4 identify $\boldsymbol{\Theta}^*$ the maximum likelihood estimation of $\boldsymbol{\Theta}$,
- 5 compute the posterior distribution of Y that is conditioned by \mathcal{F}_N , $\boldsymbol{\Theta}^*$ and \mathbf{f} .

The optimized GPR is thus given by the **mean function** of this conditioned Gaussian process, $\mathbb{E}[Y(\mathbf{x}) \mid \mathcal{F}_N, \boldsymbol{\Theta}^*, \mathbf{f}]$.

$$\mathbb{E} [Y(\mathbf{x}) \mid \mathcal{F}_N, \Theta^*, \mathbf{f}] = y^{\text{trend}}(\mathbf{x}) + y^{\text{cond}}(\mathbf{x}),$$

$$y^{\text{trend}}(\mathbf{x}) := \langle \mathbf{f}(\mathbf{x}), \hat{\beta} \rangle$$

$$y^{\text{cond}}(\mathbf{x}) := \mathbf{r}(\mathbf{x})^T [C(\Theta^*)]^{-1} (\mathbb{Y} - [F] \hat{\beta}).$$

Comments

- when N tends to infinity, the role of y^{trend} becomes negligible,
- on the contrary, when N is relatively small compared to the complexity of y , the role of y^{trend} is crucial.

- 1 Introduction
- 2 Gaussian process-based regression (GPR)
- 3 Optimization of the GPR
- 4 Nested polynomial trends
- 5 Conclusions and prospects

- When N , the number of code evaluations, is low compared to the complexity of y , such approaches are limited by the fact that **only low values** of M , the dimension of the projection family, can be considered to avoid extra-fitting.
- In order to be able to deal with higher values of M , without increasing the number of unknown parameters to be identified, we can propose a new parameterization of the polynomial trend, which is based on a **nested structure**.
- It will then be shown to what extent such an approach shows very promising results when only very limited information is available ($\leftrightarrow N$ is small).

Main ideas

- a **composition of two polynomials** for the mean function, μ , of the GP associated with $x \mapsto y(x)$ is proposed,
- the composition of polynomials being polynomials, it is thus possible to span a **large subset** of $L^2(\mathcal{D}^d, \mathbb{R})$ from a **small number** of independent parameters.

$$\mu(x; \mathbf{b}_1, \mathbf{b}_2) := \left\langle \mathbf{m}^{(p_2)} \left(\left\langle \mathbf{m}^{(p_1)}(x_1), \mathbf{b}_1 \right\rangle \right), \mathbf{b}_2 \right\rangle,$$

where $\mathbf{m}^{(p)}$ gathers all the polynomials, which total degree is less than $p \in \mathbb{N}^*$.

- $\mathbf{b}_1 \leftrightarrow$ parameterization of the inner polynomials,
- $\mathbf{b}_2 \leftrightarrow$ parameterization of the outer polynomials.

Main difficulties

- 1 the fonction

$$\mathbf{b}_1 \mapsto \left\langle \mathbf{m}^{(p_2)} \left(\left\langle \mathbf{m}^{(p_1)}(x_1), \mathbf{b}_1 \right\rangle \right), \mathbf{b}_2 \right\rangle,$$

is strongly non-linear.

- 2 different values of $(\mathbf{b}_1, \mathbf{b}_2)$ could lead to the same nested-representation.

In that prospect, it has been proposed in [Perrin et al., 2016] to work on:

- a **minimal parametrization** of this nested structure to avoid redundancies,
- a **linearization** of the mean around the maximum likelihood estimates of $(\Theta, \mathbf{b}_1, \mathbf{b}_2)$,
- **iterative algorithms** to solve the maximum likelihood maximization.

Comments on the proposed nested representation

- contrary to the "LAR+UK" approach that looks for **sparse** representations of μ , the nested representation looks for **"full"** polynomial representations, which are however characterized by a very **limited number of independent coefficients**. This is particularly interesting for the modeling of complex phenomena with very limited information.
- For $d > 1$, it allows us to model separately the **dependency** structure between the different input parameters, and the **individual** actions of each input parameter. Hence, analyzing the final structure of the mean function can give us information about the structure of y (is the model additive up to a transformation of its input parameters or are the dependencies more complex?).

For $d = 1$ and $\mathcal{D}^d = [-1, 1]$, three analytic examples are proposed:

- case 1: $y(x) = P_2^{(4)} \circ P_1^{(4)}(x)$,
- case 2: $y(x) = \sin((x + 1)^3)$,
- case 3: $y(x) = \sin(20x) \cos(2x)$,

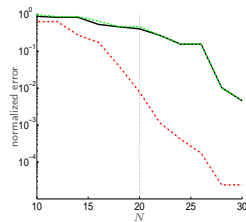
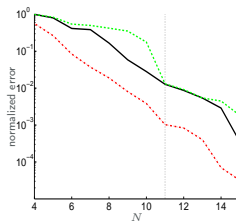
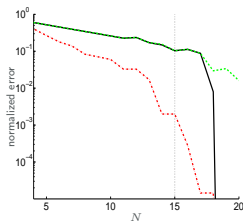


Figure: Mean value of the errors associated with these 10 repetitions. Solid black line: error associated with the LAR+UK approach. Red dotted line: error associated with the nested approach. Green dotted line: UK with linear trend

Nested polynomial trends

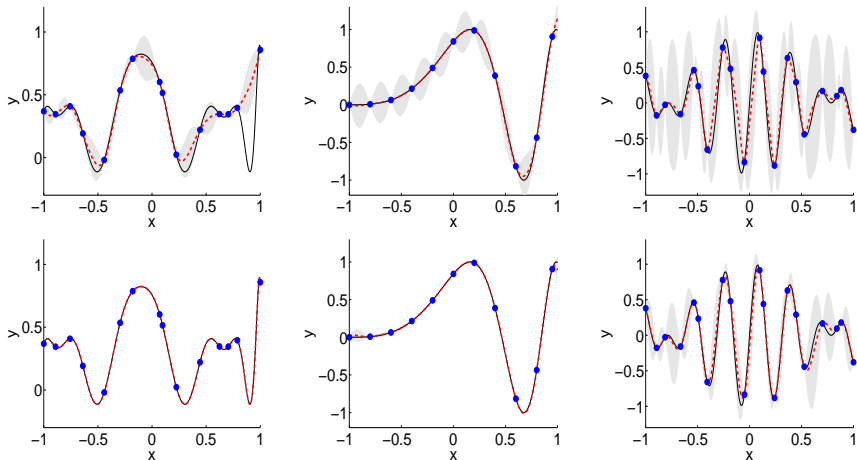


Figure: black $\leftrightarrow y_2$, blue points \leftrightarrow observations, red $\leftrightarrow \hat{\mu}_N$. top : the LAR+UK approach, bottom : the proposed approach, grey areas \leftrightarrow 95% confidence region for the prediction. Left: case 1, middle: case 2, right: case 3.

For $d > 1$, the same kinds of results are found:

- case 1: $y(\mathbf{x}) = (1 - x_1^2) \cos(7x_1) \times (1 - x_2^2) \sin(5x_2)$,
- case 2: $y(\mathbf{x}) = \sin(x_1) + 7 \sin(x_2)^2 + 0.1 \sin(x_1)x_3^4$ (**Ishigami**),
- case 3: $y(\mathbf{x}) = g^{(1)} \circ g^{(2)}(\mathbf{x})$, $g^{(1)}(\mathbf{z}) = 0.1 \cos\left(\sum_{i=1}^6 z_i\right) + \sum_{i=1}^6 z_i^2$,
 $g^{(2)}(\mathbf{x}) = (\cos(\pi x_1 + 1), \cos(\pi x_2 + 2), \dots, \cos(\pi x_6 + 6))$.

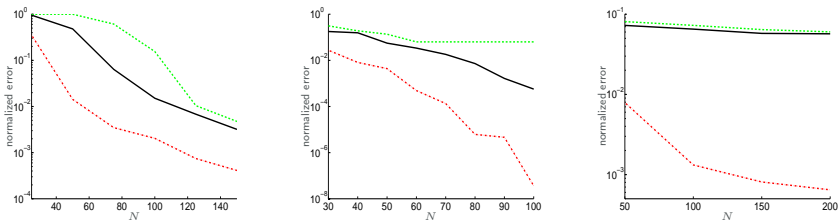


Figure: Mean value of the errors associated with these 10 repetitions. Solid black line: error associated with the LAR+UK approach. Red dashed line: error associated with the nested approach. Green dotted line: UK with linear trend.

- 1 Introduction
- 2 Gaussian process-based regression (GPR)
- 3 Optimization of the GPR
- 4 Nested polynomial trends
- 5 Conclusions and prospects

- A lot of methods in uncertainty quantification and computer experiments require **many code evaluations**.
- When the numerical cost associated with each evaluation is high, **surrogate models** are generally introduced.
- When interested in deterministic mappings (when the dimension of the input space is not too large), one of the most used method is the **GPR**.
- **Working on the polynomial trend** can strongly improve the relevance of the GPR.
- When little information about the code is available, considering **nested polynomial trends** shows promising results.
- Enabling these techniques to deal with **high dimensional problem** ($d \gg 1$), even if a lot of code evaluations ($N \gg 1$) are available is still an open questions.

Thank you for your attention.



Bachoc, F. (2013).

Estimation paramétrique de la fonction de covariance dans le modèle de Krigeage par processus Gaussiens. Application à la quantification des incertitudes en simulation numérique.

PhD thesis, University Paris Diderot, France.



Bect, J., Ginsbourger, D., Li, L., Picheny, V., and Vasquez, E. (2012).

Sequential design of computer experiments for the estimation of a probability of failure.

Statistics and Computing, 22.



Bichon, B., Eldred, M., Swiler, L., Mahadevan, S., and McFarland, J. (2008).

Efficient global reliability analysis for non linear implicit performance functions.

AIAA Journal, 46(10).



Blatman, G. and Sudret, B. (2011).
Adaptative sparse polynomial chaos expansion based on least angle regression.
Journal of Computational Physics, 230.



Dubrule, O. (1983).
Cross validation of kriging in a unique neighborhood.
Mathematical Geology, 15(6):687–699.



Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004).
Least angle regression.
Ann. Stat., 32:407–499.



Hastie, T., Tibshirani, R., and Friedman (2002).
Elements of Statistical Learning.
Springer, New York.



Kersaudy, P., Sudret, B., Varsier, N., and Picon, O. (2015).
A new surrogate modeling technique combining kriging and polynomial
chaos expansions - application to uncertainty analysis in computational
dosimetry.

Journal of Computational Physics, 286:103–117.



Miller, R. G. (1974).
The jackknife - a review.

Biometrika, 61:1–15.



Perrin, G., Soize, C., Garnier, J., and Marque-Pucheu, S. (2016).
Nested polynomial trends for the improvement of gaussian
process-based predictors.

Journal of Computational Physics, in review.



Sacks, J., Welch, W., Mitchell, T., and Wynn, H. (1989).
Design and analysis of computer experiments.
Statistical Science, 4:409–435.



Santner, T. J., Williams, B., and Notz, W. (2003).
The design and analysis of computer experiments.
Springer, New York.

