

# A Random Field Model and Decision Support in Industrial Production.

**Titre:** Un modèle par champ aléatoire et un outil d'aide à la décision en production industrielle.

Julie Oger<sup>1,2</sup>, Emmanuel Lesigne<sup>1</sup> and Philippe Leduc<sup>2</sup>

**Abstract:** We propose a new tool of decision support in front of a globally unknown phenomenon which is modeled by a random field representing simultaneously our knowledge and our lack of information. This tool is the distribution of a random variable called *failure risk probability*. Before giving the precise definition of this object, we describe an industrial context in which the decision problem occurs and we discuss Bayesian random field model constructions.

**Résumé :** Nous proposons un nouvel outil d'aide à la décision pour l'étude d'un phénomène inconnu modélisé par un champ aléatoire représentant simultanément notre connaissance et notre manque d'information. Cet outil est la distribution d'une variable aléatoire appelée *probabilité du risque de défaillance*. Avant de préciser la définition de cet objet, nous décrivons un contexte industriel dans lequel un problème décisionnel apparaît et nous examinons des constructions bayésiennes de modèles par champs aléatoires.

**Keywords:** kriging, Bayesian inference, Gaussian processes mixture prior, multivariate t-distribution, uncertainty analysis, manufacturing yield evaluation, decision support

**Mots-clés :** krigeage, inférence bayésienne, mélange de processus gaussiens, distribution de Student multivariée, analyse d'incertitude, évaluation de rendement industriel, aide à la décision

**AMS 2000 subject classifications:** 62C, 62K, 62P30

## 1. Introduction

### 1.0. Extended abstract

In competitive industries, a reliable yield forecasting is a prime factor to accurately determine the production costs and therefore ensure profitability. Indeed, quantifying the risks long before the effective manufacturing process enables fact-based decision-making. From the development stage, improvement efforts can be early identified and prioritized. In order to measure the impact of industrial process fluctuations on the product performances, the construction of a failure risk probability estimator is presented in this article. The complex relationship between the process technology and the product design (non linearities, multi-modal features...) is handled via random process regression (interpolation). A random field encodes, for each product configuration, the available information regarding the risk of non-compliance. After a brief presentation of the Gaussian model approach, we describe a Bayesian reasoning avoiding a priori choices of location

<sup>1</sup> Université François-Rabelais, CNRS, LMPT UMR7350, Tours, France.

E-mail: [emmanuel.lesigne@lmpt.univ-tours.fr](mailto:emmanuel.lesigne@lmpt.univ-tours.fr)

<sup>2</sup> STMicroelectronics, Tours, France.

E-mail: [julie.oger@cifre.org](mailto:julie.oger@cifre.org) and E-mail: [philippe.leduc@st.com](mailto:philippe.leduc@st.com)

and scale parameters. The Gaussian mixture prior, conditioned by measured (or calculated) data, yields a posterior characterized by a multivariate Student distribution.

In the second part of the article we describe the way from a random model to the failure risk probability. Our approach is to consider as random all unknown, inaccessible or fluctuating data. In order to propagate uncertainties, a fuzzy set approach provides an appropriate framework for the implementation of a Bayesian model mimicking expert elicitation. The underlying leitmotiv is to insert minimal a priori information in the failure risk model. The relevancy of this concept is illustrated with theoretical examples. Note that this article comes with supplementary material available on line (hal-00914192).

### *1.1. Settings and motivations*

During the fabrication of an industrial product, the manufacturing process cannot be entirely controlled. There is an intrinsic variability (materials, machinery, environment ...) which can significantly deteriorate the characteristics of the manufactured pieces, inducing non-functional parts. Such variations are critical when dealing with complex technical systems, such as for example integrated circuits developed in the microelectronic industry. Indeed, if this fact is covered up, the resulting circuits are (in many cases) either consistently out of specifications or unnecessarily overdesigned.

To investigate the influence of fluctuating parameters, one can perform a corner analysis applying Design of Experiments principles. However, the conclusions drawn from such measurements campaigns are, in general, qualitative and consequently incomplete. Besides, time and money required by physical prototyping are often prohibitive. As an alternative, numerical (deterministic) models implemented in engineering simulation software offer a way to compute the relevant features (thermal, mechanical, electrical ...) of a device. Thus, nowadays, engineers can virtually explore various design configurations and get a deep insight of the final product performances. In this context, statistical studies can be conducted to evaluate the effect of process tolerances on the manufacturing yield. The most popular is the Monte Carlo (MC) method, which consists in sampling configurations of the product parameters according to their probability distributions and to count the failure events. The process is monitored "in line" to roughly assess distributions of environmental variables. This method is easy to implement but its efficiency depends directly on the complexity of the deterministic model. When the numerical simulator considered is time and/or memory demanding, we only get partial information. For instance, the duration of "Finite Elements" analyses, from several hours to several days, is not compatible with a brute Monte Carlo approach.

Data are sparse and it is therefore necessary to propagate the information in the factor space using an analytical representation (emulator). Monte Carlo sampling is then applied to this analytical representation, a fast surrogate of the computer code (simulator), see [Pfingsten et al. \(2006\)](#). In their article, outputs are considered as a single outcome of a Gaussian random field. Our work focuses on a mixture of Gaussian fields, and this choice will be justified in the sequel.

Classical modeling approaches for failure rate (or yield) estimation only keep a small part of the model available information: mean, quantile... Besides, after the model has passed the validation tests, these outputs are taken at face values. We believe such a procedure is hazardous in the specific field of risk assessment. Indeed, it does not measure the impact of the uncertainty

introduced by any modeling stage on the only quantity of interest, namely the risk of failure. For example, how can the decision-maker relate the model acceptance criteria to the accuracy of the failure risk estimate? Answering this question in particular is difficult and probably fruitless. Our work proposes a novel and general solution to address these issues. Once the random model has been determined, and given the probability distribution of the product parameters, we go beyond Pfingsten's approach to define the failure risk probability. The predictive uncertainty is not deduced from the posterior model, calculating for instance the conditional variance. The failure risk itself is probabilistic and randomness mirrors the model uncertainty.

Let us be more precise and introduce the basic objects used in the sequel. We consider that the product (meaning each individual manufactured piece) under study is characterized by a number  $D$  of numerical factors. Each factor can vary in a given interval what allows the definition of the factor space  $X \subseteq \mathbb{R}^D$ . Each set of factors  $x (\in X)$  determines a numerical value  $y(x)$ , and the specifications imposed on the product apply to the value of this response  $y$ . We can derive from these specifications the *out of specifications space*  $A \subseteq \mathbb{R}$ : the product characterized by the set of factors  $x$  does not satisfy the specifications if and only if  $y(x) \in A$ . The factor space  $X$  and the out of specifications space  $A$  are considered as known, but regarding the deterministic function  $x \mapsto y(x)$ , we have only very partial information.

Besides, a probability distribution  $P$  is given on the factor space  $X$ . This distribution reflects the factors variability and is considered as known.

As we said previously, the knowledge of  $y(x)$  is rarely available for all  $x$ . We only have access to a restraint number of data. Suppose we know  $n$  deterministic response values  $(y_i = y(x_i))_{1 \leq i \leq n}$  respectively for factor set values  $(x_i)_{1 \leq i \leq n}$ . With these data, our goal is to define a random variable named *failure risk probability*. Its distribution should help for the robust estimation of the product manufacturability.

We begin with the construction of a random field model  $(Y_x)_{x \in X}$  of the unknown response  $(y(x))_{x \in X}$ .

## 1.2. Model construction

Several methods are described in the literature to infer a model  $(Y_x)_{x \in X}$  from a limited number of available deterministic data

$$Y_{x_1} = y_1, Y_{x_2} = y_2, \dots, Y_{x_n} = y_n.$$

Linear regression analysis is the method of choice in the scientific community for estimating relationships among predictor variables. When the phenomenon studied is complex though, the order of the model (usually polynomial), necessary to correctly fit the data, leads to numerous unknowns. This drawback is magnified in a high-dimensional factor space. Indeed, for a  $p$ -variable polynomial of degree  $d$ ,  $\binom{p+d}{d}$  coefficients should be determined. If the number of data available is less than this limit, the model is singular. Thus, the player is required to make a "blind bet" to enter the game... In addition, the assumptions justifying the statistical model are rarely strictly respected in practice. For example, how should we understand random errors uncorrelated with zero mean when data are obtained through computer experiments, reproducible in essence?

These difficulties are overcome by the Gaussian Process (GP) model. The original kriging method was formalized by Matheron (1989) in the geo-sciences. First, a GP indexed by  $X$  is selected and then conditioned by the data. This approach has been widely used for different applications such as, geostatistics studies (Berger et al., 2000), optimization (Emmerich et al., 2006), wind fields modeling (Cornford et al., 2002) or design sensitivity analysis (Pfingsten et al., 2006), justified by better predictive performances than several other regression methods (Rasmussen and Williams, 2006). Indeed, GP modeling has several interesting properties. Belonging to the interpolation methods family, this model exactly reproduces the observed data set, there is no residual at the observation points. As a consequence, it is an appropriate tool for the analysis of computer experiments (Sacks et al., 1989). GP can be determined (in theory) even for data set of small size, a key property when information is missing which is a common situation when the number of factors ( $D$ ) is large. Moreover, it is a very versatile model, able to describe non-continuous as well as non-differentiable surface responses. This noteworthy feature is particularly useful to handle response discontinuities, which may occur due to the numerical solving scheme of computer codes (especially when meshing algorithms are involved). Finally, the probabilistic nature of the predictions can be interpreted as a model for uncertainty, a confidence interval representing a degree of belief.

A major and often neglected problem of this method is the a priori choice of parameters (mean, variance and correlation function), left to the responsibility of the data analyst. In order to reduce the arbitrary nature of expert elicitation, we propose to replace the Gaussian field by a mixture of Gaussian fields, therefore avoiding the a priori choice of mean and variance. For this purpose, with the aim of introducing minimal a priori information into the model, mean and variance are considered as random variables with uniform probability distributions. We will show that, in doing so, any posterior random vector follows a multivariate t-distribution (or *multivariate Student distribution*). We note that a similar result has been obtained with a hierarchical Bayesian model approach described, for example, in Santner et al. (2003).

In the classical approach of yield evaluation based on a random model, the random variable

$$\omega \mapsto P(\{x \in X \mid Y_x(\omega) \in A\})$$

is considered as the model of the unknown quantity  $P(\{x \in X \mid y(x) \in A\})$  and estimated by Monte-Carlo methods. These Monte-Carlo methods can be optimized, as it is for example presented in the recent article by Auffray et al. (2014). Our approach is different, as it is explained below, and detailed in Section 3.

### 1.3. Failure risk probability viewed as a random variable

The random field  $(Y_x)_{x \in X}$  represents the information as well as the uncertainty regarding the values of the response function  $y$ . Once this field is built, we can compute for each factor point  $x \in X$ , the *failure probability*  $\mathbb{P}(Y_x \in A)$ .

The reader sees here that we denote by  $\mathbb{P}$  the probability measure associated to the random field (meaning that the random field is defined on an abstract probability space  $(\Omega, \mathbb{P})$ ). On the other hand, we recall that we denote by  $P$  the probability measure assigned to the factor space  $X$ .

We decide that there is a *failure risk* when the failure probability exceeds a fixed threshold

$\alpha \in [0, 1]$ . So, the failure risk probability is defined by  $R(\alpha) := P(\mathbb{P}(Y \in A) > \alpha)$ , that is,  $R(\alpha) = P(\{x \in X \mid \mathbb{P}(Y_x \in A) > \alpha\})$ .

Setting the correct value for the accident threshold  $\alpha$  is a tricky problem. On the one hand, the decision will significantly impact the risk assessment of the product. On the other hand, this choice is eminently subjective since it depends heavily on the risk attitude of the individual. That is why we consider  $\alpha$  as a random variable. As a consequence,  $R(\alpha)$  is also a random variable and has the distribution

$$\mathcal{R} := \int_0^1 \delta_{R(\alpha)} \eta(d\alpha)$$

where  $\eta$  is the probability distribution of  $\alpha$  and  $\delta_t$  denotes the Dirac mass at point  $t$ . We will see hereafter that the uniform distribution is a choice for  $\eta$  which provides an interesting property of mean value preservation.

In practice,  $R(\alpha)$  cannot be computed analytically, it is approximated via a Monte Carlo simulation. According to standard Bayesian methods for sampling study,  $R(\alpha)$  follows a beta distribution, denoted here by  $\beta_\alpha$ . Finally, the distribution of the failure risk probability is defined, if  $\eta$  has been chosen uniform, by

$$\mathcal{R} := \int_0^1 \beta_\alpha d\alpha.$$

This reasoning will be described with more details in the sequel.

#### 1.4. Contents

In Section 2, several constructions of random field models are discussed. After a reminder about Gaussian fields (2.1), we show that a random field prior with unknown mean (2.2) and variance leads to a conditioned random field following a multivariate t-distribution (2.3), and we discuss on the model implementation (2.4). In Section 3, the construction of the density of the failure risk probability is described, beginning with the elicitation model (3.1). We provide arguments in favor of a risk-neutral attitude, that is a uniform distribution for the accident threshold (3.2). We continue with practical considerations (3.3) to conclude with the description of our global strategy (3.4). Each section ends with illustrating examples.

## 2. From Gaussian to Student fields

The unknown function  $x \mapsto y(x)$  is modeled as a random function. In the absolute sense, if we disregard the variability that is sometimes introduced by numerical solving schemes, the function is deterministic. However, as the numerical solution to a hard problem (described, for instance, by partial differential equations),  $y$  does not have in general a closed-form expression. Consequently, the data analyst looking for a behavioral representation, can legitimately think of  $y$  as a "black-box" function: it is unknown for any particular configuration until the computer code is actually run. It hence makes sense to postulate a prior model for  $y$ , expressing our initial belief regarding  $y$  features. Bayesian updating then combines the evidences acquired, the data output  $(y_i = y(x_i))_{1 \leq i \leq n}$ , and the prior distribution to yield the posterior distribution.

Contrary to what is often implicitly stated in the literature, we believe that it is far from obvious to get a relevant prior representation for  $y$ . Standard priors reflect more of their mathematical

tractability than a real understanding of the phenomenon under study. Incorporating doubtful prior information into the model may yield erroneous and overly confident forecasts, a dangerous cocktail in risk assessment. For that reason, we focus in the following on a prior based on quite few assumptions about how the model output relates to its inputs; it is deliberately weakly informative.

### 2.1. Gaussian field

In this first sub-section, we point out that the conditioning of a Gaussian field generates a new Gaussian field and we recall classical formulae giving the conditional densities.

A Gaussian field  $(Y_x)_{x \in X}$  is characterized by the fact that all finite dimensional marginal distributions are Gaussian and by the following data:

- mean values:  $\mu(x) := \mathbb{E}(Y_x)$
- covariances:  $\rho(x, x') := \text{cov}(Y_x, Y_{x'})$ .

**Theorem 1.** *Let  $x_1, x_2, \dots, x_n \in X$  and  $y_1, y_2, \dots, y_n \in \mathbb{R}$ .*

*The conditional distribution of the Gaussian field  $(Y_x)_{x \in X}$  given  $(Y_{x_i} = y_i)_{1 \leq i \leq n}$  is still Gaussian with*

- mean values:

$$\mathbb{E}(Y_x | (Y_{x_i} = y_i)_{1 \leq i \leq n}) = \mu(x) + \rho(x, (x_i)) \Sigma^{-1} (y - \mu((x_i)))^T$$

- covariances:

$$\text{cov}(Y_x, Y_{x'} | (Y_{x_i} = y_i)_{1 \leq i \leq n}) = \rho(x, x') - \rho(x, (x_i)) \Sigma^{-1} \rho(x', (x_i))^T$$

where  $\rho(x, (x_i)) := (\rho(x, x_1), \rho(x, x_2), \dots, \rho(x, x_n))$ ,  $\Sigma := (\rho(x_i, x_j))_{1 \leq i, j \leq n}$  is a positive-definite matrix,  $y := (y_1, y_2, \dots, y_n)$  and  $\mu((x_i)) := (\mu(x_1), \mu(x_2), \dots, \mu(x_n))$ .

This is a classical result in Probability Theory.

A standard solution to our initial problem is to use this theorem with a given a priori Gaussian field  $(Y_x)_{x \in X}$  where neither the mean value  $\mu = \mu(x)$  nor the variance  $\sigma^2 = \rho(x, x)$  depend on  $x$ . In this case, we will denote  $k(x, x') = \rho(x, x') / \sigma^2$ , the correlation coefficient between  $Y_x$  and  $Y_{x'}$ . If the covariance depends only on the difference between  $x$  and  $x'$ , the Gaussian field is said stationary or homogeneous (Abrahamsen, 1997). Moreover, if the covariance only depends on the Euclidean distance between  $x$  and  $x'$ , the Gaussian field is said isotropic.

In practice, the mean  $\mu$ , the variance  $\sigma^2$  and the correlation function  $k$  are usually a priori defined by the expert or estimated by means of calibration methods. The most commonly used is the maximum likelihood method which will be discussed in Section 2.4.2.

Note that here and in the sequel, we avoid particular choices of the correlation  $k$ , vectors  $x_1, x_2, \dots, x_n$  and numbers  $y_1, y_2, \dots, y_n$  which could lead to degenerate distributions.

### 2.2. Gaussian mixture with random mean

In this sub-section, the prior model is a mixture of Gaussian processes, which differ by translation. This way, it is not required to set a priori the mean value. The posterior distribution of this parameter is deduced from the data via Bayesian inference. We will see that a (non informative)

improper prior for the mean parameter can be defined as the limit of a (weakly informative) proper uniform distribution. The limiting step is totally justified here since, on the one hand, it leads to proper posteriors and, on the other hand, inferences have a positive miscalibration (see Gelman (2006) for details), that is an overestimate (on average) of the variance. The conditioned random field is Gaussian and we give formulae for the mean and variance. Doing so, we recognize expressions initially established by Sacks et al. (1989) and also derived by Santner et al. (2003).

We propose to consider a prior random field  $Y_x = U + W_x$  where  $U$  is a real random variable following a uniform distribution on an interval  $[-m, m]$  and where  $(W_x)_{x \in X}$  is a centered Gaussian field with constant variance. Moreover, we suppose that  $U$  and  $(W_x)_{x \in X}$  are independent.

The parameters characterizing the Gaussian field  $(W_x)_{x \in X}$  are the variance  $\sigma^2$  and the correlation function  $k$  (recall that the mean is zero). So, for all  $x, x' \in X$ ,  $\text{cov}(W_x, W_{x'}) = \sigma^2 k(x, x')$  and, for all  $x \in X$ ,  $k(x, x) = 1$ .

Let  $x_1, x_2, \dots, x_n \in X$  and  $y := (y_1, y_2, \dots, y_n) \in \mathbb{R}^n$ . Denote  $\Sigma := (k(x_i, x_j))_{1 \leq i, j \leq n}$  the positive-definite matrix of correlations and  $k(x) := (k(x, x_j))_{1 \leq j \leq n}$  the correlation vector.

**Theorem 2.** *The conditional distribution of the random field  $(Y_x)_{x \in X}$  knowing that  $(Y_{x_i} = y_i)_{1 \leq i \leq n}$  is given by explicit formulae for the densities of finite dimensional marginals.*

*When the parameter  $m$  goes to infinity, this conditional distribution becomes Gaussian. In particular, when  $m \rightarrow \infty$ , the univariate conditional distribution of the random variable  $Y_x$  becomes Gaussian with mean*

$$\mu + k(x)\Sigma^{-1}(y - \mu\mathbf{1})^T \quad \text{with} \quad \mu := \frac{y\Sigma^{-1}\mathbf{1}^T}{\mathbf{1}\Sigma^{-1}\mathbf{1}^T} \quad (1)$$

and variance

$$\sigma^2 \left( 1 - k(x)\Sigma^{-1}k(x)^T + \frac{\left(1 - \mathbf{1}\Sigma^{-1}k(x)^T\right)^2}{\mathbf{1}\Sigma^{-1}\mathbf{1}^T} \right) \quad (2)$$

where  $\mathbf{1} = (1, 1, \dots, 1) \in \mathbb{R}^n$ .

*Remark.* The mean  $\mu + k(x)\Sigma^{-1}(y - \mu\mathbf{1})^T$  can also be written

$$\left( \frac{\mathbf{1}\Sigma^{-1}}{\mathbf{1}\Sigma^{-1}\mathbf{1}^T} (1 - k(x)\Sigma^{-1}\mathbf{1}^T) + k(x)\Sigma^{-1} \right) y^T.$$

Note that an expression similar to (1) is obtained in Section 2.5.

*Proof of Theorem 2.* We look for the distribution of the random field  $(Y_x)_{x \in X}$  given  $(Y_{x_i} = y_i)_{1 \leq i \leq n}$ . Let  $r$  be a positive integer and  $(t_1, \dots, t_r) \in X^r$ . We have:

$$(Y_{t_1}, \dots, Y_{t_r}, Y_{x_1}, \dots, Y_{x_n}) = (U, \dots, U) + (W_{t_1}, \dots, W_{t_r}, W_{x_1}, \dots, W_{x_n})$$

and  $(W_{t_1}, \dots, W_{t_r}, W_{x_1}, \dots, W_{x_n})$  follows the distribution  $\mathcal{N}(0, \Delta)$  with

$$\Delta := \sigma^2 \begin{pmatrix} \Sigma_2 & k(t) \\ k(t)^T & \Sigma \end{pmatrix}, \text{ a positive-definite matrix}$$

where  $k(t) := (k(t_i, x_j))_{1 \leq i \leq r, 1 \leq j \leq n}$  and  $\Sigma_2 := (k(t_i, t_j))_{1 \leq i \leq r, 1 \leq j \leq r}$ .  
Denote by  $f$  the density of the random vector  $(Y_{t_1}, \dots, Y_{t_r}, Y_{x_1}, \dots, Y_{x_n})$ :

$$f(\zeta) = \frac{1}{2^m} \int_{-m}^m \frac{1}{(\sqrt{2\pi})^{n+r} \sqrt{|\Delta|}} \exp\left(-\frac{1}{2}(\zeta - u)\Delta^{-1}(\zeta - u)^T\right) du$$

where  $\zeta := (y_{t_1}, y_{t_2}, \dots, y_{t_r}, y_{x_1}, y_{x_2}, \dots, y_{x_n})$  and  $u := (u, u, \dots, u) \in \mathbb{R}^{n+r}$ .

The conditional density of  $(Y_{t_1}, Y_{t_2}, \dots, Y_{t_r})$  given  $(Y_{x_i} = y_{x_i})_{1 \leq i \leq n}$  is

$$z := (y_{t_1}, y_{t_2}, \dots, y_{t_r}) \mapsto \frac{f(z, y)}{\int_{\mathbb{R}^r} f(z', y) dz'} =: g(z).$$

After simplification, we get:

$$g(z) = \frac{\int_{-m}^m \exp\left(-\frac{1}{2}((z \ y) - u)\Delta^{-1}((z \ y) - u)^T\right) du}{\int_{\mathbb{R}^r} \int_{-m}^m \exp\left(-\frac{1}{2}((z' \ y) - u)\Delta^{-1}((z' \ y) - u)^T\right) dz' du}.$$

By monotone convergence, some calculations show that:

$$g(z) \xrightarrow{m \rightarrow +\infty} \frac{\exp\left(-\frac{1}{2}\left((z \ y)\Delta^{-1}(z \ y)^T - \frac{(\mathbf{1}\Delta^{-1}(z \ y)^T)^2}{\mathbf{1}\Delta^{-1}\mathbf{1}^T}\right)\right)}{\int_{\mathbb{R}^r} \exp\left(-\frac{1}{2}\left((z' \ y)\Delta^{-1}(z' \ y)^T - \frac{(\mathbf{1}\Delta^{-1}(z' \ y)^T)^2}{\mathbf{1}\Delta^{-1}\mathbf{1}^T}\right)\right) dz'}.$$

Within the exponential in the numerator of this expression we identify a non-negative second degree polynomial of the variable  $z$ . We recognize a Gaussian distribution (of dimension  $r$ ).

In the limit situation  $m \rightarrow \infty$ , we note that the distribution of the conditioned random field is well-defined and is Gaussian. Some calculations show that, at the point  $x \in X$ , the one-dimensional marginal Gaussian distribution of the field has mean

$$\mu + k(x)\Sigma^{-1}(y - \mu\mathbf{1})^T \quad \text{with} \quad \mu := \frac{y\Sigma^{-1}\mathbf{1}^T}{\mathbf{1}\Sigma^{-1}\mathbf{1}^T}$$

and variance

$$\sigma^2 \left(1 - k(x)\Sigma^{-1}k(x)^T + \frac{(1 - \mathbf{1}\Sigma^{-1}k(x)^T)^2}{\mathbf{1}\Sigma^{-1}\mathbf{1}^T}\right).$$

□

### 2.3. Gaussian mixture with random mean and variance

We go one step further, avoiding the choice of the mean and variance. The prior model is a mixture of Gaussian processes, which differ by affine transformation. Once again, the conditioned random

field is well-defined selecting (weakly informative) proper uniform priors for the unknown mean and standard deviation. When the distributions supports respectively tend to the entire and positive real lines, the multivariate Student distribution arises as the random field posterior distribution. We give explicit formulae for the location and scale parameters. Santner et al. (2003) and Kato (2009) describe similar models with different priors. For example, with Jeffreys priors for the mean and standard deviation (scale-invariance property), Santner and al. come to the same posterior except for the number of degrees of freedom. We justify our choice arguing that a positive miscalibration is always preferred by the risk-averse individual.

### 2.3.1. Multivariate Student distribution

First, recall the definition of a multivariate Student distribution. We refer to the book of Kotz and Nadarajah (2004). A  $p$ -variate Student distribution (or a multivariate t-distribution) has density:

$$t \mapsto \frac{1}{(\sqrt{\pi v})^p \sqrt{|\Sigma|}} \frac{\Gamma\left(\frac{v+p}{2}\right)}{\Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{1}{v} (t - \mu) \Sigma^{-1} (t - \mu)^T\right)^{-\frac{v+p}{2}} \quad (t \in \mathbb{R}^p)$$

where the positive integer  $v$  is the number of degrees of freedom,  $\Sigma$  is the  $p \times p$  positive-definite matrix of scale parameters and  $\mu$  is the  $1 \times p$  vector of location parameters.

For  $v > 1$ , the mean vector of the Student distribution is well-defined and equals to  $\mu$ . For  $v > 2$ , the covariance matrix of the Student distribution is well-defined and equals to  $\frac{v}{v-2} \Sigma$ .

This is a multi-dimensional generalization of the Student distribution. When  $v = 1$ , the distribution is a multivariate Cauchy distribution. When  $v$  goes to infinity, the distribution tends to a multivariate Gaussian distribution.

### 2.3.2. Student field posterior

Now, we propose to consider a prior random field  $Y_x = U + VW_x$ , where  $U$  is a real random variable following a uniform distribution on an interval  $[-m, m]$ , where  $V$  is a real and positive random variable following a uniform distribution on an interval  $[\varepsilon, 1/\varepsilon]$  and where  $(W_x)_{x \in X}$  is a centered normalized Gaussian field. Moreover, we suppose that  $U$ ,  $V$  and  $(W_x)_{x \in X}$  are independent.

The parameter characterizing the Gaussian field  $(W_x)_{x \in X}$  is the correlation function  $k$  (recall that the mean is zero and the variance is 1). We suppose here that  $n \geq 3$ . Let  $x_1, x_2, \dots, x_n \in X$  and  $y := (y_1, y_2, \dots, y_n) \in \mathbb{R}^n$ . Denote  $\Sigma := (k(x_i, x_j))_{1 \leq i, j \leq n}$  the positive-definite matrix of correlations and  $k(x) := (k(x, x_j))_{1 \leq j \leq n}$  the correlation vector.

**Theorem 3.** *The conditional distribution of the random field  $(Y_x)_{x \in X}$  knowing that  $(Y_{x_i} = y_i)_{1 \leq i \leq n}$  is given by explicit formulae of densities of finite dimensional marginals.*

*When the parameter  $m$  goes to infinity and  $\varepsilon$  goes to zero, for  $n > 2$ , this conditional distribution becomes a multivariate Student distribution.*

*In particular, when  $m \rightarrow \infty$ ,  $\varepsilon \rightarrow 0$  and  $n > 2$ , the univariate conditional distribution of the random variable  $Y_x$  becomes a Student distribution with  $n - 2$  degrees of freedom, with location parameter*

$$\mu + k(x) \Sigma^{-1} (y - \mu \mathbf{1})^T \quad \text{with} \quad \mu := \frac{y \Sigma^{-1} \mathbf{1}^T}{\mathbf{1} \Sigma^{-1} \mathbf{1}^T}$$

and scale parameter

$$\sqrt{\frac{1}{n-2} ((y - \mu \mathbf{1}) \Sigma^{-1} y^T) \left( 1 - k(x) \Sigma^{-1} k(x)^T + \frac{(1 - \mathbf{1} \Sigma^{-1} k(x)^T)^2}{\mathbf{1} \Sigma^{-1} \mathbf{1}^T} \right)}$$

where  $\mathbf{1} = (1, 1, \dots, 1) \in \mathbb{R}^n$ .

The proof of Theorem 3 results from the same reasoning than the proof of Theorem 2 with more complex formal calculations. It is developed in the supplementary materials available on line (hal-00914192).

*Remark.* The expression  $(y - \mu \mathbf{1}) \Sigma^{-1} y^T$  can also be written  $(y - \mu \mathbf{1}) \Sigma^{-1} (y - \mu \mathbf{1})^T$ . Note that a similar expression is obtained in Section 2.5.

## 2.4. Model implementation

### 2.4.1. Discussion on correlation functions

The a priori choice of the correlation function is one of the major difficulties in random field modeling (Sacks et al., 1989) and the discussion is unfortunately often avoided. It is sometimes justified by expert knowledge (Cornford et al., 2002) in particular applied cases. Rasmussen and Williams (2006) give a complete list of correlation functions, classified according to practical considerations. The most popular correlation families are isotropic: Exponential, Matern and Rational Quadratic classes (Abrahamsen, 1997). Indeed, an isotropic random field does not suffer the curse of dimensionality since the number of parameters in the model does not depend on the factor space dimension ( $D$ ).

Determining appropriate values for the parameters of the correlation function is the purpose of calibration methods. The most popular of them is the maximum likelihood estimation (MLE) described in the next sub-section. Rasmussen and Williams (2006), Stein (1999) and Robert (2007) describe this method and discuss about its capabilities and limitations such as numerical issues for too big samples, multiple optima for too small samples and over-fitting problems for which the sample is well learned but the unknown function values are poorly predicted everywhere else in the factor space.

In order to overtake these different limitations, a good way seems to try different choices of correlation function; each alternative is then evaluated by checking if known values of the responses are compatible with the confidence intervals predicted by the random field. The literature describes several types of cross-validation methods measuring the predictive capability of a model (Currin et al., 1988). It is of course necessary to use a set of data which is not involved in the construction of the a posteriori model. Because we cannot, when information is scarce, afford not to include all the gathered data in the final model, the partition into training and validation sets is only temporary, for the purpose of the cross-validation stage. In that case, the test is therefore performed on an "incomplete" posterior model version.

Rasmussen and Williams (2006) describe an alternative method which consists in choosing a parameterized family of correlation function and a prior distribution on its parameters, to construct

a hierarchical model. This method involves analytical approximations of integrals. Markov chain Monte Carlo (MCMC) methods are popular solutions to make these computations, see [Robert \(2007\)](#) for detailed description of these methods. A major drawback is here the cost of such calculations.

#### 2.4.2. Maximum Likelihood Estimation

The classical maximum likelihood estimation (MLE) method defines an estimator of the unknown parameter vector  $\theta$  of a probability distribution  $f_\theta$ . This estimator is the value  $\theta_{\max}$  which maximizes the density distribution of a random sample calculated at the observed value of this sample.

Under the same name, this method has been adapted to the identification of the parameters of an unknown random field  $(Y_x)_{x \in X}$  when considering a family of values  $Y_{x_1}, Y_{x_2}, \dots, Y_{x_n}$ . It is important to have in mind that the observation set is a single outcome of the random field so that the validity of the MLE method, in this context, is not obvious. However, this topic is outside the scope of this article.

Considering now the probabilistic modeling of a deterministic phenomenon, the MLE provides a practical procedure to set the parameters of the a priori random field  $(Y_x)_{x \in X}$  knowing  $Y_{x_i} = y_i$ ,  $1 \leq i \leq n$ . This is a classical approach in Gaussian field modeling, see [Rasmussen and Williams \(2006\)](#), [Stein \(1999\)](#) and [Robert \(2007\)](#).

Even if our justifications are incomplete, let us describe the MLE for the random model described in Section 2.3.

First, set a parametric family of correlation functions  $k_\theta$  depending on the parameter vector  $\theta$ , which defines also the parametric correlation matrix  $\Sigma_\theta = (k_\theta(x_i, x_j))_{1 \leq i, j \leq n}$ . The MLE consists in choosing  $\theta$  which maximizes the density of the random vector  $(Y_{x_1}, Y_{x_2}, \dots, Y_{x_n})$  according to  $m$  and  $\varepsilon$ :

$$f_{m,\varepsilon}(\theta) = \frac{1}{2m} \frac{1}{\frac{1}{\varepsilon} - \varepsilon} \int_{-m}^m \int_{\varepsilon}^{\frac{1}{\varepsilon}} \frac{1}{(\sqrt{2\pi})^n v^n \sqrt{|\Sigma_\theta|}} \exp\left(-\frac{1}{2v^2} (y-u)\Sigma_\theta^{-1}(y-u)^T\right) du dv$$

with  $y := (y_{x_1}, y_{x_2}, \dots, y_{x_n})$  and  $u := (u, u, \dots, u) \in \mathbb{R}^n$ .

The maximum likelihood estimator  $\theta_{\max}(m, \varepsilon)$  maximizes  $f_{m,\varepsilon}(\theta)$ . It maximizes also:

$$\tilde{f}_{m,\varepsilon}(\theta) = \int_{-m}^m \int_{\varepsilon}^{\frac{1}{\varepsilon}} \frac{1}{v^n \sqrt{|\Sigma_\theta|}} \exp\left(-\frac{1}{2v^2} (y-u)\Sigma_\theta^{-1}(y-u)^T\right) du dv.$$

We don't know any analytic expression of  $\tilde{f}_{m,\varepsilon}(\theta)$  and a fortiori of  $\theta_{\max}(m, \varepsilon)$  but since we are interested in large values of  $m$  and little values of  $\varepsilon$ , we consider the limit value  $\tilde{f}_{\infty,0}(\theta)$ . So we propose to study:

$$\tilde{f}_{\infty,0}(\theta) = \int_{-\infty}^{+\infty} \int_0^{+\infty} \frac{1}{v^n \sqrt{|\Sigma_\theta|}} \exp\left(-\frac{1}{2v^2} (y-u)\Sigma_\theta^{-1}(y-u)^T\right) du dv.$$

A short calculation gives

$$\tilde{f}_{\infty,0}(\theta) = \frac{\sqrt{\pi} 2^{\frac{n-3}{2}} (n\sigma_\theta^2)^{\frac{2-n}{2}} \Gamma\left(\frac{n-2}{2}\right)}{\sqrt{|\Sigma_\theta|} (\mathbf{1}\Sigma_\theta^{-1}\mathbf{1}^T)}$$

with

$$\sigma_{\theta}^2 := \frac{1}{n} (y - \mu_{\theta} \mathbf{1}) \Sigma_{\theta}^{-1} (y - \mu_{\theta} \mathbf{1})^T \quad \text{and} \quad \mu_{\theta} := \frac{\mathbf{1} \Sigma_{\theta}^{-1} y^T}{\mathbf{1} \Sigma_{\theta}^{-1} \mathbf{1}^T}.$$

Maximizing  $\tilde{f}_{\infty,0}(\theta)$  is equivalent to minimize  $-\ln(\tilde{f}_{\infty,0}(\theta))$  and after simplifications, we get:

$$\theta_{\max}(\infty, 0) \quad \text{minimizes} \quad n \ln(\sigma_{\theta}^2) + \ln(|\Sigma_{\theta}|) + 2 \ln\left(\frac{\mathbf{1} \Sigma_{\theta}^{-1} \mathbf{1}^T}{\sigma_{\theta}^2}\right).$$

Observe the difference with the Gaussian case in which the MLE method proposes analytical expressions for the mean and the variance of the field and where the parameter vector  $\theta$  of the correlation function is estimated by minimizing:

$$n \ln(\sigma_{\theta}^2) + \ln(|\Sigma_{\theta}|).$$

### 2.5. Comparison between Student and Gaussian fields

Let us come back for a while to the Gaussian field model described in Section 2.1, choosing constant mean and variance. If the mean  $\mu$  and variance  $\sigma^2$  are determined by the MLE method considering the data set  $(y(x_i) = y_i)_{1 \leq i \leq n}$ , then the following formulae are obtained:

$$\mu = \frac{y \Sigma^{-1} \mathbf{1}^T}{\mathbf{1} \Sigma^{-1} \mathbf{1}^T} \quad \text{and} \quad \sigma^2 = \frac{1}{n} (y - \mu \mathbf{1}) \Sigma^{-1} (y - \mu \mathbf{1})^T.$$

This is a classical result in this area. See for example [Currin et al. \(1988\)](#) or [Stein \(1999\)](#).

Following Theorem 1, the a posteriori mean and variance of the Gaussian random variable  $Y_x$  are respectively

$$\mu + k(x) \Sigma^{-1} (y - \mu \mathbf{1})^T \quad \text{and} \quad \sigma^2 \left(1 - k(x) \Sigma^{-1} k(x)^T\right).$$

It is interesting to compare these values with those given in Theorem 3. Note that, for the same correlation function, the mean remains the same and the variance is greater. We can illustrate this difference on a simple example in one dimension.

Let  $X = [-5, 5]$ . Let a sample of 5 points  $(x_1 = -4, x_2 = -3, x_3 = -1, x_4 = 0, x_5 = 2)$  with the associated response values  $(y_1 = -2, y_2 = 0, y_3 = 1, y_4 = 2, y_5 = -1)$ . We build for an arbitrary fixed correlation function,  $\rho(x, x') := \exp(-100|x - x'|^2)$ , a Gaussian field and a Student field given the knowledge of the previous sample.

Figure 1 shows, for each field after conditioning, the mean and the boundaries of an interval computed with a confidence of 0.90. As expected, the confidence intervals are both null at the 5 sampled points since the distributions at these points degenerate in Dirac distributions. The confidence interval of the Student posterior is larger than the Gaussian posterior equivalent, suggesting a lower degree of belief in the prediction. The Student field seems more reliable, since we take into account all the possible values of the mean  $\mu$  and the variance  $\sigma^2$ , that were arbitrarily fixed in the case of the Gaussian field.

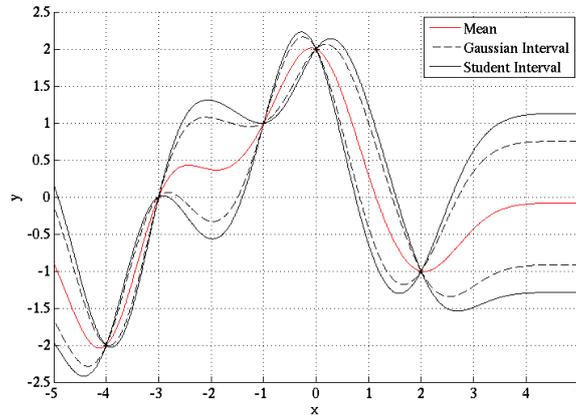


FIGURE 1. Univariate illustration. Gaussian and Student field posteriors derived from five evaluations of the unknown function  $y(x)$ .

### 3. Decision support

We present here a new approach in order to define a failure risk probability. As it is briefly explained in Introduction, it applies to any deterministic phenomena on which we have only partial information, when this information and our uncertainty are described by a random field model.

#### 3.1. Failure risk probability

We use here notations introduced at the end of Sub-section 1.1. The *real failure risk* is the measure of the *real failure set*  $\{x \in X \mid y(x) \in A\}$  where  $X \subseteq \mathbb{R}^D$  is the factor space provided with a probability distribution  $P$ . In a lot of practical cases, the function  $x \mapsto y(x)$  is not available though; it can only be sampled at a few  $x$  locations. According to Section 2, a random field  $(Y_x)_{x \in X}$  defined on a probability space  $(\Omega, \mathbb{P})$  is derived from a set of observations in order to approximate the function  $y$ .

This modeling stage is not without consequences on the approach retained to calculate the failure risk probability. Since the knowledge on the function  $y$  is low, it seems reasonable to provide the failure risk prediction with a confidence measure. Thus a decision-maker could easily evaluate the model quality and rule on the manufactured product robustness more safely. This problem can be formulated as follows: how can we propagate the uncertainty inherent to any prediction based on a random field to the failure risk estimation?

If the phenomenon under study was really a random field  $(Y_x)_{x \in X}$ , we could consider the random variable

$$Z := \omega \mapsto P(\{x \in X \mid Y_x(\omega) \in A\}) \quad \text{on } (\Omega, \mathbb{P})$$

as the *failure risk* and its distribution as the *failure risk probability*. This is the definition usually retained in risk assessment studies. See for example the work of Auffray et al. (2014).

Strictly speaking, the model  $(Y_x)_{x \in X}$  cannot describe the reality, which is perfectly determined (but unknown). It correctly represents our knowledge at the observation points. Besides, we have to give an interpretation to the process randomness. If we retain the subjective conception of probability theory, randomness results from the incomplete knowledge of the quantity  $y(x)$ . Indeed, several realizations of the conditioned random field meets the requirements of the modeling problem, that is interpolates the observations. Within this framework, the measure  $\mathbb{P}$  quantifies a degree of belief regarding the model forecasts. Randomness becomes a way to assess the quality, i.e. the predictive capability of the model. The expression  $\mathbb{P}(Y_x \in A) = 0.5$  does not mean that the point  $x$  has a one-in-two chance of belonging to the real failure set; it means that, regarding the available information, one random field realization out of two states that  $x$  belongs to the real failure set. The model is therefore inadequate at point  $x$ .

Even if the previous definition of the failure risk by the random variable  $Z$  is absolutely correct mathematically, we claim that it cannot be used, in this particular case, to model incomplete knowledge. Indeed, the distribution of  $Z$  can be a Dirac mass even in the case of a bad model. An example is necessary here to clarify our point. According to the Lemma demonstrated in Appendix A.1, if the measure  $P$  is continuous and if the random variables  $Y_x$  are pairwise independent and satisfy  $\mathbb{P}(Y_x \in A) = \varepsilon$ , then the distribution of  $Z$  is a Dirac mass in  $\varepsilon$ . As a consequence, the particular case where  $\varepsilon$  is near 1/2, which is the sign of a poor model, is not incompatible with the fact that the distribution of the random variable  $Z$  is a Dirac mass<sup>1</sup>. However, a Dirac mass distribution, or a distribution close to a Dirac mass, represents a phenomenon for which the information is sure, or close to sure. This does not reflect the degree of insecurity associated to an inadequate model. Consequently, we conclude that the distribution of the random variable  $Z = P(Y(\omega) \in A)$  does not provide relevant information regarding the model quality.

The random field alone is not a knowledge representation; it cannot encode the mental procedures involved in belief assessments. An additional analysis is hence necessary. We certainly do not have at our disposal a strict condition of membership to the real failure set. Nevertheless, we can easily compute the following *failure probability*, a random variable seen as a membership function  $\mathcal{M}$ :

$$\begin{cases} (X, P) \rightarrow [0, 1] \\ x \mapsto \mathcal{M}(x) := \mathbb{P}(Y_x \in A). \end{cases}$$

The quantity  $\mathcal{M}(x)$  is the proportion of the random field outcomes concluding that  $x$  belongs to the failure set. Such a function, only involves marginal distributions, and can be straightforwardly evaluated. It associates to each point  $x \in X$  a real number in the interval  $[0, 1]$  measuring a grade of membership of  $x$  to the failure set. Thus, the crisp failure set, unknown, can be only imprecisely characterized by the membership function  $\mathcal{M}$ . Due to incomplete knowledge, we can only access to the fuzzy version of the failure set. According to Zadeh (1965), "such a framework provides a natural way of dealing with problems in which the sources of imprecision is the absence of sharply defined criteria of class membership rather than the presence of random variables."

By definition, the notion of "belonging" is not well-defined for a fuzzy set. For a given grade of membership  $\mathcal{M}(x)$ , a decision-maker may consider that  $x$  belongs to the failure set whereas another may not, depending on its own risk tolerance. In order to account for the subjective nature of the decision, we introduce the level  $\alpha$ ,  $0 \leq \alpha \leq 1$ , and agree to say that " $x$  belong the failure

<sup>1</sup> See Appendix A.2 for a detailed description and discussion of this toy example.

set" if  $\mathcal{M}(x) > \alpha$ . The threshold  $\alpha$  determines the status to give to uncertain predictions. These data are therefore aggregated in two classes: "good" and "poor". Let us define the  $\alpha$ -level failure set:

$$\{x \in X \mid \mathcal{M}(x) > \alpha\}.$$

It is a crisp subset of  $X$  and can be consequently measured:

$$R(\alpha) := P(\{x \in X \mid \mathcal{M}(x) > \alpha\}).$$

It is the failure risk probability estimation by a decision-maker whose risk tolerance is  $\alpha$ . We implicitly define here a causal model: the risk probability is conditional on the risk tolerance  $\alpha$  (Pearl, 1988). The decision-maker risk tolerance being unknown, we choose to consider the threshold  $\alpha$  as random. It takes its values in  $[0, 1]$  and we denote by  $\eta$  its distribution. In this framework, the *failure risk probability* is finally defined as a random variable whose distribution is the image of the measure  $\eta$  by the function  $R$ . In other words, this distribution is

$$\mathcal{R} := \int_0^1 \delta_{R(\alpha)} \eta(d\alpha).$$

### 3.2. Uniform distribution for the tolerance threshold

In this sub-section, we discuss the choice of the distribution  $\eta$  of the threshold  $\alpha$ .

A natural output in risk assessment is the average failure probability,  $E(\mathbb{P}(Y \in A))$  in our case. It is an estimator of the  $P$ -measure of the real failure set, based on knowledge of the values  $(y_i = y(x_i))_{1 \leq i \leq n}$ . As we will see in the sequel, the causal model introduced in Section 3.1 can be interpreted as an operator acting on the space of probability measures on  $[0, 1]$ , which links the distributions of the random variables  $R$  (failure risk probability) and  $\mathcal{M}$  (failure probability). It seems clever to keep the first moment of the latter when the transform is applied. In that way, the mean value of the failure risk probability remains an estimator of the failure set measure. This remark leads us to the first proposition.

If the distribution of the failure probability is uniform, we cannot conclude about the model relevance, defined as its ability to discriminate the factor space (i.e. identify the failure set). In such situation, the distribution of the failure risk probability should also be uniform, meaning that the state of uncertainty has been preserved. Based on this desired invariance property, we state the second proposition.

First of all, we introduce some definitions and notations, useful for the subsequent demonstrations. As above, for each  $\alpha \in [0, 1]$ , we set:

$$R(\alpha) := P(\mathcal{M}(x) > \alpha).$$

The function  $R$  is the complementary cumulative distribution function (or tail distribution) of the random variable  $\mathcal{M}(x)$ , defined on the probability space  $(X, P)$ . Now, we can also look at  $R$  as a random variable defined on the probability space  $([0, 1], \eta)$ , and we define, for all  $t \in [0, 1]$ :

$$G(t) := \eta(\{\alpha \in [0, 1] \mid R(\alpha) > t\}) = \eta(R > t).$$

So  $G$  is the tail distribution of the random variable  $R$ .

Moreover, we denote by  $K$  the cumulative distribution function of  $\eta$ :  $K(u) = \eta([0, u])$  for all  $u \in [0, 1]$ . Finally, let us recall the definition of the generalized inverse of the decreasing function  $R$ .

**Definition** (Generalized inverse). Let  $S$  a decreasing right continuous function on the interval  $[0, 1]$  with  $S(0) \leq 1$  and  $S(1) = 0$ . For all  $t \in [0, 1]$ , we set:

$$S^{-1}(t) = \sup \{ \alpha \in [0, 1] : S(\alpha) > t \}.$$

At this stage, it is not difficult to check that  $G = K \circ R^{-1}$ . Let us define in a general setting the operator which maps  $G$  onto  $R$ . The probability measure  $\eta$  is assumed fixed. To any probability measure  $m$  on the interval  $[0, 1]$ , we can associate another probability measure  $m'$  on  $[0, 1]$  in the following way:

- firstly denote by  $F_m$  the function  $F_m(\alpha) = m((\alpha, 1])$ ,
- secondly define  $F_{m'} = K \circ F_m^{-1}$ .

Denote by  $L_\eta$  the operator, acting on the space of probability measures on  $[0, 1]$ , which associates  $m'$  to  $m$ .

**Proposition 1.**

1. The operator  $L_\eta$  preserves the first moment of the probability (i.e.  $\int_0^1 t m(dt) = \int_0^1 t m'(dt)$  for all  $m$ ) if and only if the probability  $\eta$  is uniform on  $[0, 1]$ .
2. If the probability  $\eta$  is uniform, then the operator  $L_\eta$  is its own inverse, meaning that  $(m')' = m$ .

*Proof of Proposition 1.* The first moment of the probability measure  $m$  associated with  $F_m$  is equal to:

$$\int_0^1 F_m(\alpha) d\alpha.$$

The following neat identity can be proved:

$$\int_0^1 F_m(\alpha) d\alpha = \int_0^1 F_m^{-1}(\alpha) d\alpha.$$

Thus, the first moment of  $m$  is equal to the first moment of  $m'$  if and only if:

$$\int_0^1 K(F_m^{-1}(\alpha)) d\alpha = \int_0^1 F_m^{-1}(\alpha) d\alpha.$$

It can be easily verified that the identity map is the unique function  $K$  which satisfies this equality for all choice of  $m$ . This means that  $\eta$  has to be a uniform distribution. It is not difficult to check point 2. □

**Proposition 2.** *The uniform distribution  $\eta$  is the unique distribution for which the operator  $L_\eta$  applied to a uniform distribution gives a uniform distribution.*

The proof does not present any difficulty.

In the sequel,  $\eta$  is chosen as the uniform distribution. It results from proposition 1 that  $E_\eta(R) = E(\mathbb{P}(Y \in A))$ , where  $E(\mathbb{P}(Y \in A))$  is an estimator of the real failure risk  $P(y \in A)$ .

At this stage of our presentation, the distribution of the failure risk probability is defined as follow. The relevance of this formula is illustrated in Appendix A.2, on the toy example introduced in Section 3.1.

$$\mathcal{R} := \int_0^1 \delta_{R(\alpha)} d\alpha.$$

### 3.3. Monte Carlo approach

We describe in this sub-section why we have to add a new layer of Bayesian randomness, due to our practical approach of the calculation of the risk probability  $R(\alpha)$ .

In practice, we use a MC method to get a numerical approximation of  $R(\alpha)$ . More precisely, an importance sampling is performed: independent samples  $(x_m)_{1 \leq m \leq M}$ , in the factor space  $X$ , are generated from the distribution  $P$ . For more details on the related preferential MC method, see 4.5 p. 19 of [Caffish \(1998\)](#). For each point  $x_m$ , the failure probability  $\mathcal{M}(x_m) = \mathbb{P}(Y_{x_m} \in A)$  can be computed numerically since the Student distribution is tabulated. Denote by  $n(\alpha)$  the number of MC draws  $x_m$  such that  $\mathcal{M}(x_m) > \alpha$ . The MC estimator of  $R(\alpha)$  is the quotient  $n(\alpha)/M$ .

Note that a classical Bayesian reasoning is used to model the uncertainty due to finite sampling. In this approach, the MC sampling is considered as a binomial experiment. This assumption is thoroughly justified for a good quality pseudo random number generator. As a consequence,  $R(\alpha)$  is a random variable whose distribution is the beta distribution  $\beta_\alpha$  with shape parameters  $(n(\alpha) + 1, M - n(\alpha) + 1)$ , assuming a uniform prior. Adding this additional probabilistic stratum in the model, we may appear a little bit pernickety. Nevertheless, it guards us against the basic (but widespread) fallacy which consists in presuming that a system is perfectly safe as long as no accident has been observed.

In conclusion, taking into account the numerical aspects of the problem, we define the distribution of the failure risk probability as

$$\mathcal{R} := \int_0^1 \beta_\alpha d\alpha.$$

### 3.4. Global strategy

In this sub-section, we gather together Sections 2 and 3, and we describe the global strategy that has been implemented in the practical application of our present work. From a set of virtual experiments, the computer code emulator (Student field) is derived. A stochastic simulation is then performed (uncertainty analysis) to extract the distribution of the failure risk probability. Standard statistical quantities such as the mean, the standard deviation and a confidence interval are available to the decision-maker for risk assessment. The confidence interval indicates the reliability of the estimate, the sources of uncertainty being the partial knowledge of the original function  $y$  and the finite MC sampling.

If the prediction quality is too low, new data points from the space  $X$  have to be added to the observations set  $\{x_1, x_2, \dots, x_n\}$ . Let us describe succinctly a process of choice for the new data points. Note first that the MC simulation only requires model evaluations. Consequently, we have full scope to complete the observations set; it is in no way linked to the factors distribution  $P$ . We could sample data uniformly in  $X$  but, as each datum "costs the earth", it is strongly

recommended to structure the data collection process. This issue is addressed by experimental design techniques, which intend to optimize the information gathering. We retain the differential entropy  $h(Y_x) := \mathbb{E}(-\log f(Y_x))$ , where  $f$  is the density function of the random variable  $Y_x$ , as a measure of the lack of information (Cover and Thomas, 2006). The information provided to the probabilistic model is increased if observations are made at locations in  $X$  where the entropy is maximum. We face a multimodal optimization problem : the objective is to find a set of local maxima for  $h(Y_x)$ . It can be solved by conventional methods such as gradient descent or metaheuristics as evolutionary algorithms. We refer to Talbi (2009) for a complete presentation of metaheuristics. Once the new observation points are identified, we are back to step one... The iterative procedure is repeated until the precision requirements regarding the failure risk probability are met or the due date for the risk analysis is reached.

### 3.5. Example and results

In order to illustrate the interest of the risk assessment scheme presented in this article, three theoretical examples described in Appendix B are considered. Using test functions with closed-form expressions, the failure probability can be calculated exactly. Therefore, estimates can be compared to the true value.

The Gaussian mixture prior defined in Section 2.3 is kitted out with the anisotropic  $\gamma$ -exponential correlation function (see Appendix C) and trained using the MLE method described in Section 2.4.2. The resulting model has  $D + 1$  parameters, a reasonable level of complexity. The factor distribution  $P$  is assumed uniform.

As an alternative to our model-based Monte Carlo (MMC) method, we consider the brute-force Monte Carlo (BMC) method plus the Bayesian model for finite sampling, recalled in Section 3.3. Let  $k$  the number of defectives samples, i.e. such that  $y(x) \in A$ , from  $M$  trials. The inferred failure risk probability follows the beta distribution with shape parameters  $(k + 1, M - k + 1)$ . Note that this trivial model implicitly assumes a sampling of the factor space according to the distribution  $P$  (uniform). We used a quasi-random rather a pseudo-random source:  $(x_m)_{1 \leq m \leq M}$  are chosen as elements of the Sobol low discrepancy sequence. Such sampling covers the factor space more evenly, a desirable property to obtain a relevant statistical population. For the MMC and the BMC methods, the mean value and the boundaries of the smallest confidence interval, stated at the 90% confidence level, are plotted as a function of the number  $M$  of evaluations,  $10 \leq M \leq 1000$ .

Before going further, it is important to have in mind that a confidence interval only provides a statistical estimation of the error on the result. It obviously does not imply a strict condition of membership to the interval. Indeed, whatever the sampling method used, this one may not catch the essential features of the function  $y$ . As a result, uncertainty can be underestimated, a dangerous situation in risk assessment. This is particularly true for BMC, which does not take into account the spatial distribution of the data (geometrical structure of the space  $X$ ). The third example well illustrates this point. This is the key difference with the MMC method we propose. Besides, and this is the cornerstone of this article, we aim at taking into account the uncertainty introduced by the interpolation stage. As a consequence, the comparison of the convergence properties of MMC and BMC should not only focus on confidence intervals. The uncertainty model is more relevant in the second case.

### 3.5.1. Quadric example

The characteristic of the MMC method is to warn the user if the emulator is poor. It is justified in practice, though, if it can outperform the BMC method at least in some cases.

So, let us start with the quadric example, considering a space  $X$  of dimension  $D = 5$  and a true failure risk  $r_t = 1.008 \times 10^{-1}$  (the theoretical expression is available in the supplementary materials). See Figure 2 for a comparison of the convergence rates of the model-based and brute-force Monte Carlo. It turns out that, in this case, MMC definitely outperforms BMC since it estimates the real failure risk as accurate as BMC on 600 samples, requiring a third of the function evaluations.

Although MMC has an extra source of uncertainty, its confidence interval is smaller. At least 60 samples are necessary to get a confidence interval size inferior to 0.05. This lower bound soars to 400 data points using BMC. Besides, the MMC confidence interval always includes the failure risk true value. All this suggests that the prior assumptions of the model are appropriate and significantly impact the prediction accuracy, which is not surprising, the function under study being quite simple.

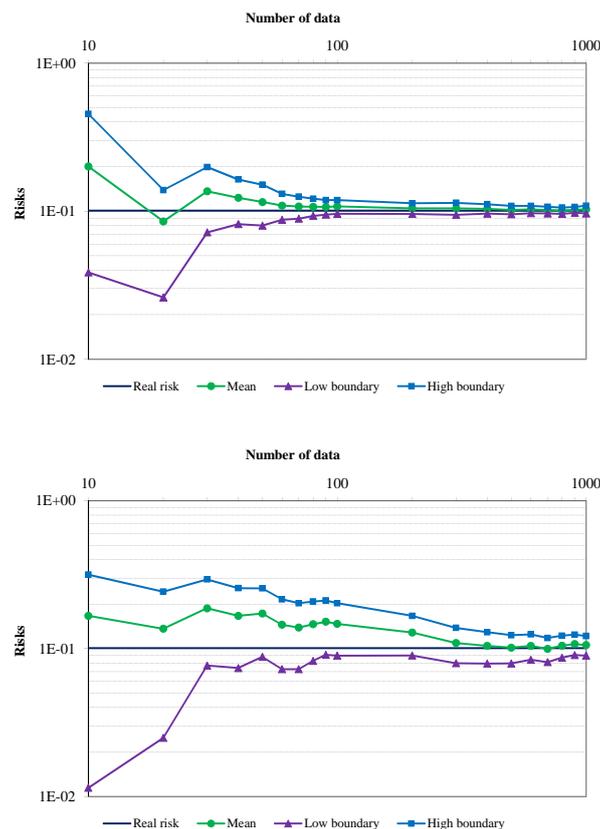


FIGURE 2. *Quadric example: convergence plots for MMC method (top) and BMC method (bottom)*

### 3.5.2. Sine example

The sine example ( $D = 3$ ) is a trickier problem since the function to be modeled is oscillating. However, we expect the BMC method to perform correctly since the true failure risk is set to  $r_t = 2.048 \times 10^{-1}$ , a relatively high value.

Both methods give quite similar results, plotted in Figure 3. The relative error on the mean failure risk predicted by MMC is less than 5% beyond 300 samples, to be compared to 400 samples for BMC. The MMC and BMC confidence intervals size is less than 0.05 respectively beyond 600 and 800 samples.

We note that, below 200 samples, the MMC confidence interval is large, which indicates that the conditioning of the prior model does not restrict much the set of interpolating functions. The propagation of information is therefore weak and we conclude that the  $\gamma$ -exponential correlation function does not really suit. In addition, a spurious pinch of the confidence interval shows up for 30 samples, due to strong variations of the correlation function parameters. The maximization of the marginal likelihood yields an overly confident estimate of the correlation function parameters posterior that would be obtained by classical Bayesian analysis. Thus, MLE reaches its limits when information is scarce, firstly because it becomes very sensitive to the training set and secondly because the Bayesian posterior can be hardly approximated by its mode.

### 3.5.3. Bell-shaped example

Consider the bell-shaped example ( $D = 5$ ). The benchmark function, as a mixture of 10 multi-dimensional Gaussian functions, is hilly. The failure set is disconnected: it is the union of 10 small hyperspheres of different radii. The real failure risk is  $r_t = 4.341 \times 10^{-3}$ .

Below 700 draws, the statistical mean decreases as  $O(n^{-1})$ , which is typical of an under-sampling of the failure set. Indeed, the first defective part has been observed beyond 800 draws. At the end of the sampling process ( $M = 1000$ ), the relative error on the mean failure risk is still 23% for MMC and 31% for BMC. BMC is here obviously overconfident: the upper boundary of the 90% confidence interval comes close to touching the risk true value. MMC prediction is much more robust: the confidence interval safely flanks  $r_t$  and the ratio of the mean value to the confidence interval size unambiguously indicates that additional data are required.

The raw analysis of the prediction error of  $y$  gives no evidence to conclude that information is lacking. Indeed, the average relative error, introduced by the approximation  $\mathbb{E}(Y_x)$  of  $y$ , is about 1.7%, which is not exactly the signature of a poor model. In contrast with global sensitivity analysis, discussed in Oakley and O'Hagan (2004), goodness-of-fit requirements do not only depend upon the visiting probability distribution  $P$ : good accuracy is also required in this case at the failure set boundaries. Thus, the decision regarding the model adequacy should not rest on a global criterion. In addition, focusing on the model performances to validate the failure risk prediction may lead to over-quality: the accident set  $A$  should also be taken into account. The probabilistic risk assessment scheme we propose naturally overcome these difficulties.

## 4. Conclusion

We introduce in this article a data sparing scheme in order to measure the robustness of a design to fabrication fluctuations. It was evolved in order to offset the lack of information, either because of

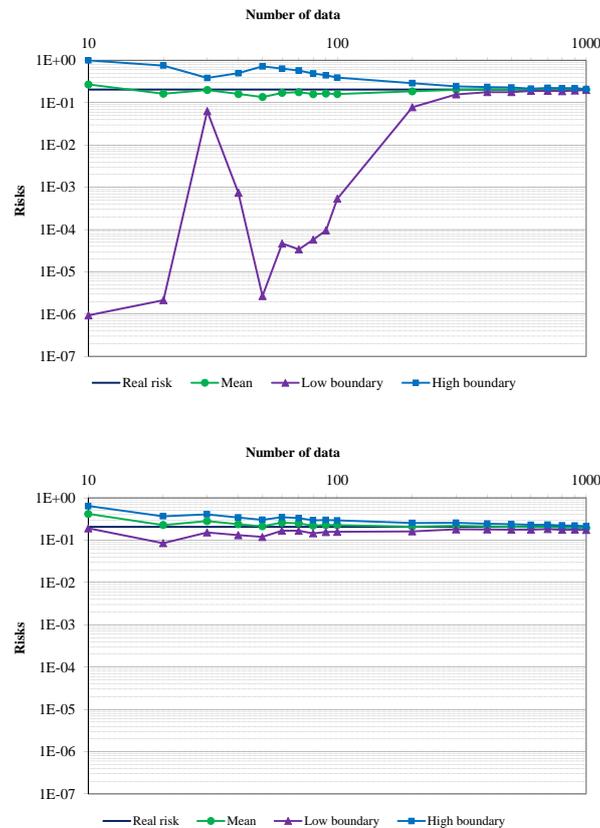


FIGURE 3. Sine example: convergence plots for MMC method (top) and BMC method (bottom)

the complexity of the data acquisition process or due to the dimensionality of the state space. The method is therefore particularly adapted to virtual experimentation, in which physically realistic computer experiments replace test runs to foresee product performances. Indeed, simulations duration are often prohibitive; the numerical simulator is considered as a "black-box" function and approximated by an analytical emulator. In this work, the regression is based on a Student process derived from low-informative priors (uniform) on the location and scale parameters of a Gaussian process.

Risk evaluation can be viewed as a binary classification problem, the product specifications defining the class membership. Thus, we use the behavioral model as a probabilistic classifier, by interpreting predictive probabilities as degrees of belief. In this scope, the failure risk probability naturally comes up as a random variable, randomness modeling the attitude of the decision-maker while exposed to uncertainty. As a result, the distribution of the failure risk probability is a straightforward measure of the impact of the manufacturing process variability on the product performances, in a way that reflects the various sources of uncertainty: incomplete knowledge of the "black-box" function, values of the model variables, numerical approximations...

The relevance of the proposed methodology has been demonstrated on theoretical examples.

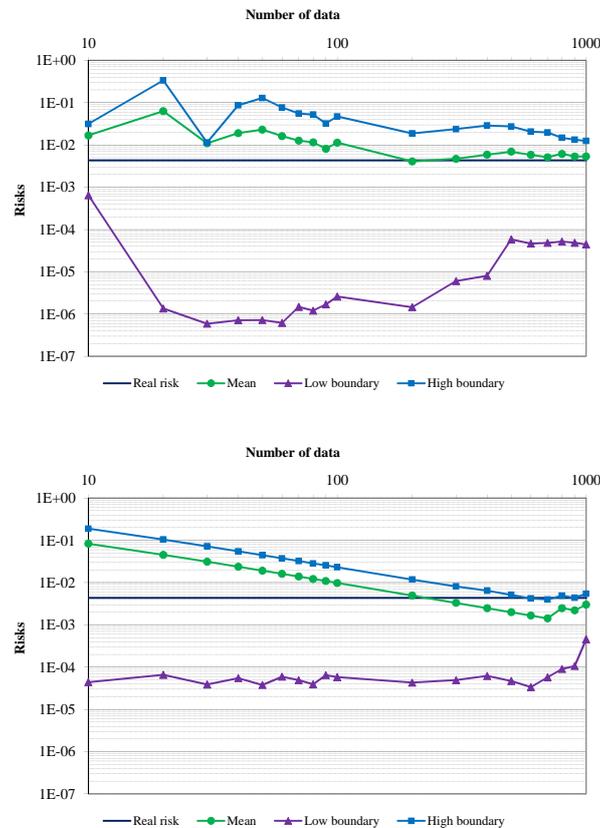


FIGURE 4. *Bell-shaped example: convergence plots for MMC method (top) and BMC method (bottom)*

This study shows that, for risk assessment, model-based Monte Carlo provides a more reliable estimator of the failure rate than brute-force Monte Carlo. Although it is not a crucial point in our argument, we noticed that introducing an intermediary analytical representation (i.e. incorporating prior knowledge) may come with a higher convergence rate of the MC sampling, thus reducing the computational effort for a given accuracy.

This work opens up new perspectives for the application of random field regression to engineering risk analysis and we can already give further lines of research:

- A natural extension of this work is the generalization of the approach in order to address the multi-responses problem. A solution is to build a random field over a mixed state space including the factors as well as the responses. The main drawback here is that the space dimension may be increased drastically. In addition, it may be difficult to find an appropriate a priori correlation function. A possibly clever alternative is to model each response independently and to compute the upper bound of the failure probability of the union event. The "worst case" risk probability is then easily derived, applying the procedure prescribed in this article.

- We can study a more sophisticated regression model of the type

$$\sum_{j=1}^k U_j f_j(x) + VW$$

(where  $f_j$  are deterministic known functions). It is the Student extension of some classical Gaussian models extensively described in the literature (see Santner et al., 2003).

- The following ternary classification problem can be investigated. If the MC samples are separated into three groups: "good", "bad" and "undecidable", it is possible to define a deficiency risk. Such random variable is a quantitative measure of the random model inability to detect defective parts. It is intended to notify explicitly the decision-maker that the model is not informed enough to draw reliable conclusions.
- In practice, we also have to control carefully numerical calculations and approximations. A particularly resistant problem is the definite-positive matrix inversion when the matrix is closed to singular.

## Acknowledgments

This research has been developed through a PhD thesis, supported by STMicroelectronics and the *Association Nationale de la Recherche et de la Technologie* (<http://www.anrt.asso.fr/>) via a CIFRE funding.

The authors thank the anonymous referee of this article for stimulating questions which enable them to clarify some points of the presentation.

## Appendix A

### A.1. Lemma

**Lemma.** *If the probability  $P$  on the set  $X$  is continuous, the random variables  $Y_x$  are pairwise independent and satisfy  $\mathbb{P}(Y_x \in A) = \varepsilon$  for all  $x$ , then  $P(Y_x(\omega) \in A) = \varepsilon$  for almost all  $\omega$ .*

*Proof of the lemma.*

$$\begin{aligned} \mathbb{E} \left[ (P(Y \in A))^2 \right] &= \int_{\Omega} \left[ \int_X \mathbf{1}_A(Y_x(\omega)) dP(x) \cdot \int_X \mathbf{1}_A(Y_{x'}(\omega)) dP(x') \right] d\mathbb{P}(\omega) \\ &= \int_X \int_X \left[ \int_{\Omega} \mathbf{1}_A(Y_x(\omega)) \cdot \mathbf{1}_A(Y_{x'}(\omega)) d\mathbb{P}(\omega) \right] dP(x) dP(x') \end{aligned}$$

From the fact that  $P$  is continuous, we deduce that for  $P \otimes P$  almost all  $(x, x')$ , we have  $x \neq x'$ , hence  $Y_x$  and  $Y_{x'}$  are independent random variables. We obtain

$$\begin{aligned} \mathbb{E} \left[ (P(Y \in A))^2 \right] &= \int_X \int_X \left[ \int_{\Omega} \mathbf{1}_A(Y_x(\omega)) \mathbb{P}(\omega) \cdot \int_{\Omega} \mathbf{1}_A(Y_{x'}(\omega)) d\mathbb{P}(\omega) \right] dP(x) dP(x') \\ &= \int_X \int_{\Omega} \mathbf{1}_A(Y_x(\omega)) dP(x) d\mathbb{P}(\omega) \cdot \int_X \int_{\Omega} \mathbf{1}_A(Y_{x'}(\omega)) dP(x') d\mathbb{P}(\omega), \end{aligned}$$

which gives

$$\mathbb{E} \left[ (P(Y \in A))^2 \right] = [\mathbb{E}(P(Y \in A))]^2.$$

As a consequence the random variable  $P(Y \in A)$  is deterministic, its distribution is a Dirac mass.

Since we have

$$\mathbb{E}[P(Y \in A)] = E[\mathbb{P}(Y \in A)] = \varepsilon,$$

we conclude that the distribution of the random variable  $P(Y \in A)$  is the Dirac mass at  $\varepsilon$ .  $\square$

## A.2. A persuasive test case

The toy example, introduced in the section 3.1, is now detailed and discussed in order to answer these two frequently asked questions :

1. *Why the distribution of the random variable  $Z := \omega \mapsto P(\{x \in X \mid Y_x(\omega) \in A\})$  is not the object that contains the information we need ?*

In order to show that the desirable information is not contained in the distribution of  $Z$ , let us give an example. Let  $\varepsilon$  be any number between 0 and 1. We want to describe two extreme situations (somewhat caricatural, but it is possible to approximate them by realistic situations). In these two radically different situations, the distribution of  $Z$  will be a Dirac mass at  $\varepsilon$ .

- (i) First is the case where the random field  $(Y_x)$  gives an exact (or excellent) description of the function  $x \mapsto y(x)$ , and where  $P(y \in A) = \varepsilon$ . In this case, we have a perfect (excellent) model and, for all  $x \in X$ ,  $\mathbb{P}(Y_x \in A) = 0$  or 1.
- (ii) Second is an example of bad model. Consider the case where the probability  $P$  on the set  $X$  is continuous, the random variables  $Y_x$  are pairwise independent and satisfy  $P(Y_x \in A) = \varepsilon$  for all  $x$ . Under these conditions, we have  $P(Y(\omega) \in A) = \varepsilon$  for almost all  $\omega$ . (In Appendix A.1, we give the proof of this claim.)

2. *What is the information given by the random variable called failure risk probability and denoted by  $R$  ?*

A necessary condition for a good model is the fact that the distribution of  $R$  is concentrated. Coming back to our example above, we see that in case (i) the distribution of  $R$  is a Dirac mass at  $\varepsilon$ , whereas in case (ii) the distribution of  $R$  is  $(1 - \varepsilon)\delta_0 + \varepsilon\delta_1$ . (We denote here by  $\delta_a$  the Dirac mass at the point  $a$ ).

## Appendix B: Examples

Discussions in the Section 3.5 are based on three arbitrary chosen benchmark functions that we describe briefly here and more precisely in supplementary material (hal-00914192) available on line. In each of these examples, the out of specification space is of the form  $A := [m, +\infty)$  for some  $m \in \mathbb{R}$ .

**B.1. Quadric example**

Fix  $m > 0$ ,  $a_1, a_2, \dots, a_D > m^2$ ,  $X := [-1, 1]^D$  and for  $x := (x_1, x_2, \dots, x_D) \in X$ :

$$y(x) := \sqrt{\sum_{i=1}^D a_i x_i^2}.$$

**B.2. Sine example**

Fix  $a_1, a_2, \dots, a_D$  non zero integers,  $X := [0, 1]^D$  and for  $x := (x_1, x_2, \dots, x_D) \in X$ :

$$y(x) := \sin \left( 2\pi \sum_{i=1}^D a_i x_i \right).$$

**B.3. Bell-shaped example**

Let  $R$  be a positive integer,  $X = [0, 1]^D$  and for  $x \in X$ :

$$y(x) := \max_{1 \leq i \leq R} f_i(x, \mu_i, \sigma_i)$$

where for all  $1 \leq i \leq R$ ,  $\mu_i \in \mathbb{R}^D$ ,  $\sigma_i \in \mathbb{R}$ ,  $\sigma_i > 0$  and

$$f_i(x, \mu_i, \sigma_i) = \frac{1}{(\sqrt{2\pi}\sigma_i)^D} \exp \left( -\frac{1}{2\sigma_i^2} \|x - \mu_i\|^2 \right)$$

such that for the fixed threshold  $m$ , the failure areas associated to each of the  $f_i$  do not overlap each other.

**Appendix C: Correlation Function**

The correlation function chosen in the Section 3.5 belongs to  $\gamma$ -exponential family:

$$\rho(x, x') := \exp \left( - \left( \sum_{i=1}^D \left( \frac{x_i - x'_i}{l_i} \right)^2 \right)^{\frac{\gamma}{2}} \right) \quad \text{for } x, x' \in X.$$

It is stationary and anisotropic. Roughly speaking, the exponent  $\gamma \in (0, 2]$  controls the smoothness of the random process (it is mean square differentiable only when  $\gamma = 2$ ) and the positive real numbers  $l_1, l_2, \dots, l_D$  are characteristic length-scales, defining the influence hyper-ellipsoid of an observation point.

## References

- Abrahamsen, P. (1997). *A review of Gaussian random fields and correlation functions - 2nd edition*. Norwegian Computing Center.
- Auffray, Y., Barbillon, P., and Marin, J.-M. (2014). Bounding rare event probabilities in computer experiments. *Computational Statistics and Data Analysis*, 80(0):153–166.
- Berger, B., De Oliviera, V., and Sansó, B. (2000). Objective bayesian analysis of spatially correlated data. *Journal of the American Statistical Association*, 96:1361–1374.
- Caffish, R. (1998). Monte carlo and quasi-monte carlo methods. *Acta Numerica*, 7:1–49.
- Cornford, D., Nabney, I., and Williams, C. (2002). Modelling frontal discontinuities in wind fields. *Journal of nonparametric statistics*, 14(1-2):43–58.
- Cover, T. and Thomas, J. (2006). *Elements of information theory, 2nd edition*. Wiley-Interscience.
- Currin, C., Mitchell, T., Morris, M., and Ylvisaker, D. (1988). A bayesian approach to the design and analysis of computer experiments. Technical report, Technical Report ORNL-6498.
- Emmerich, M., Giannakoglou, K., and Naujoks, B. (2006). Single- and multiobjective evolutionary optimization assisted by gaussian random field metamodels. *IEEE Transactions on Evolutionary Computation*, 10:421–439.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, 1:515–533.
- Kato, K. (2009). Improved prediction for a multivariate normal distribution with unknown mean and variance. *Annals of the Institute of Statistical Mathematics*, 61:531–542.
- Kotz, S. and Nadarajah, S. (2004). *Multivariate t-distributions and their applications*. Cambridge University Press.
- Matheron, G. (1989). *Estimating and choosing. An essay on probability in practice. Translated from the French and with a preface by A. M. Hasofer*. Springer-Verlag.
- Oakley, J. and O’Hagan, A. (2004). Probabilistic sensitivity analysis of complex models: a bayesian approach. *Journal of the Royal Statistical Society*, 66:751–769.
- Pearl, J. (1988). Do we need higher-order probabilities, and if so, what do they mean? In *Conference on Uncertainty in Artificial Intelligence*.
- Pfingsten, T., Herrmann, D., and Rasmussen, C. (2006). Model-based design analysis and yield optimization. *IEEE Transactions on semiconductor manufacturing*, 19:475–486.
- Rasmussen, C. and Williams, C. (2006). *Gaussian processes for machine learning*. The MIT Press.
- Robert, C. (2007). *The Bayesian choice: from decision-theoretic foundations to computational implementation - 2nd edition*. Springer-Verlag New York Inc.
- Sacks, J., Welch, W., Mitchell, T., and Wynn, H. (1989). Design and analysis of computer experiments. *Statistical Science*, 4:409–435.
- Santner, T., Williams, B., and Notz, W. (2003). *The design and analysis of computer experiments*. Springer.
- Stein, M. (1999). *Interpolation of spatial data: some theory for kriging*. Springer.
- Talbi, E. (2009). *Metaheuristics: from design to implementation*, volume 74. Wiley.
- Zadeh, L. (1965). Fuzzy sets. *Information and control*, 8:338–353.