

AN ALTERNATIVE ESTIMATOR OF A PROBABILITY OF FAILURE

LUCIE BERNARD*

*Univ. Pierre et Marie Curie, Laboratoire LSTA
75005 Paris, France
email: lucie.bernard@upmc.fr*

PHILIPPE LEDUC

*STMicroelectronics
37000 Tours, France
email: philippe.leduc@st.com*

We are interested in estimating a failure probability, defined as a threshold-exceeding probability of a random variable which is costly to simulate and whose distribution is unknown. In view of the restricted number of available observations of this random variable, classical Monte Carlo simulation methods can not be used. That is why we adopt a Bayesian approach and assume that the failure probability is a realization of a random variable relied on the so-called Gaussian process regression model. In order to provide a reliable estimation of the failure probability, it is desirable to learn as much as possible about the posterior distribution of this random variable. Considering that this is not obvious, we propose an alternative random variable whose good properties, in term of simulation, improve the estimation of the failure probability. In particular, we show that there exists a convex order between these two variables and we exploit the resulting properties.

Keywords: Failure probability, Bayesian methods, Gaussian random process, Convex order.

1. Context

During the fabrication of an industrial product, the manufacturing process cannot be entirely controlled. Some operations are complicated to manage and inevitably subjected to variability. For instance, succession of heavy and complex machinery or fabrication areas whose temperature and humidity are challenging to maintain constant over time contain tasks which are difficult to reproduce identically. As a result, it is inevitable that some

*Corresponding Author

finished products do not fulfil the imposed specifications, which may have negative effects on their performance and more generally on the profit. In competitive industries, a reliable yield forecasting is then a prime factor to accurately determine the production costs and therefore ensure profitability. Our goal is to measure, long before the manufacturing process be effective, the impact of manufacturing process variations on the production profitability through the estimate of a failure probability.

Our theoretical framework is the following: considering the inherent variability, the product under study has $d \geq 1$ design parameters, usually called factors, whose numerical values vary in a given interval. This allows to consider the definition of the factor space $\mathbb{X} \subseteq \mathbb{R}^d$ and the introduction of a random variable \mathbf{X} defined over the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, taking values in the measurable space $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$. The joint probability distribution $P_{\mathbf{X}}$ of \mathbf{X} is known. We say that a product is characterized by \mathbf{x} when its factors take the values $\mathbf{X} = \mathbf{x}$.

The performance of a product characterized by \mathbf{x} is measured from the output of a numerical simulation (also called computer experiment), which consists in evaluating a measurable and deterministic function $g : \mathbb{X} \rightarrow \mathbb{R}$ at point \mathbf{x} . More precisely, if the output $g(\mathbf{x})$ exceeds a prescribed threshold $T \in \mathbb{R}$, then the product characterized by \mathbf{x} does not satisfy the imposed specifications and is considered as non-functional. Thus, in order to make a reliable estimation of the profitability, we are interested in approximating the threshold-exceeding probability

$$p = P_{\mathbf{X}}(\{\mathbf{x} \in \mathbb{X} : g(\mathbf{x}) \geq T\}) = \mathbb{P}(g(\mathbf{X}) \geq T) = \int_{\mathbb{X}} \mathbb{1}_{g(\mathbf{x}) \geq T} P_{\mathbf{X}}(d\mathbf{x}), \quad (1)$$

typically called the probability of failure.

Here, the function g has no available analytical expression (g is a black-box function) and is expensive to evaluate. As a result, we cannot use a crude Monte Carlo estimator to provide an estimate of the failure probability. We only observe evaluations of the function g at points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. We denote by \mathbf{g}^n the vector of observations $(g(\mathbf{x}_1), \dots, g(\mathbf{x}_n))$.

2. Gaussian process regression

In this context, a natural idea is to adopt a Bayesian point of view and consider the observations \mathbf{g}^n as being incomplete information about a realization of a random process ξ indexed by \mathbb{X} .

Let $\boldsymbol{\xi}^n$ be the vector $(\xi(\mathbf{x}_1), \dots, \xi(\mathbf{x}_n))$. By choosing a Gaussian a priori probability distribution for the process ξ , it is, conditionally to $\boldsymbol{\xi}^n = \mathbf{g}^n$, still Gaussian. In the sequel, in order to simplify the notations, we denote by ξ_n the posterior process, i.e. the process ξ conditionally to $\boldsymbol{\xi}^n = \mathbf{g}^n$.

Note that in the Gaussian case, each realization of the process ξ_n , i.e. each function $\mathbf{x} \mapsto \xi_n(\mathbf{x}, \cdot)$, goes through the points $(\mathbf{x}_i, g(\mathbf{x}_i))$, $i = 1, \dots, n$. In other words, the realizations of the process ξ_n are interpolation functions of the observed values of the function g . See for instance Figure 1 and Figure 2, where the prior process ξ and the posterior process ξ_n are respectively represented.

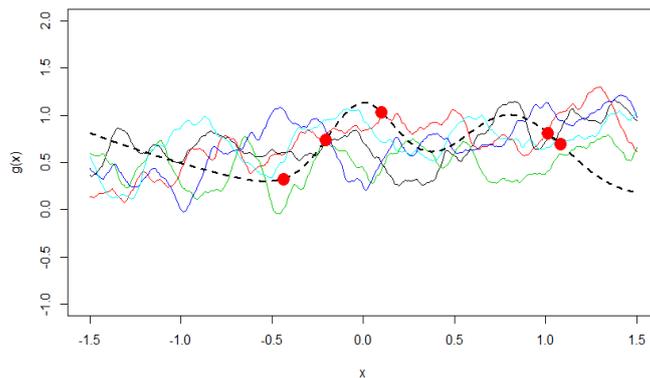


Figure 1. In dimension $d = 1$. The red points are the observations $(\mathbf{x}_i, g(\mathbf{x}_i))_{1 \leq i \leq n}$. The black dashed is the true function g . Other lines are realizations of the Gaussian process ξ .

This approach refers to the so-called Gaussian process regression method or *Kriging*. As mentioned in [5] and [6], a Gaussian distribution as a prior probability distribution for ξ offers in our context better predictive performance than several other regression methods. Besides, under the Gaussian assumption, the computations required for inference and learning become relatively easy. Indeed, the mean and variance function of the process ξ_n are explicitly known (see e.g. [8]).

Here, it is not necessary to explicitly know the distribution of the posterior process. The results presented further are valid for any choice of a prior

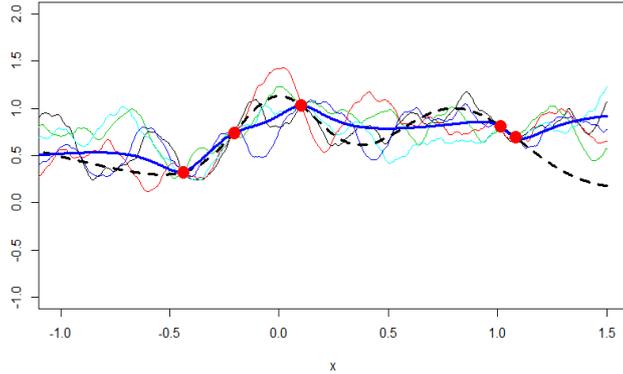


Figure 2. In dimension $d = 1$. The black dashed line is the true function g . The red points are the observations $(\mathbf{x}_i, g(\mathbf{x}_i))_{1 \leq i \leq n}$. The full blue line is the mean function of the Gaussian process ξ_n . Other lines are realizations of the Gaussian process ξ_n .

distribution of ξ . However, we made the assumption that it is Gaussian in order to make the calculations easier.

3. A first estimation of the probability of failure

Assuming that the function g is a realization of a Gaussian process, it seems natural to consider that the probability p a realization of the random variable $P : \Omega \rightarrow [0, 1]$ defined by:

$$P = \mathbb{P}(\xi(\mathbf{X}) \geq T \mid \xi) = \int_{\mathbb{X}} \mathbb{1}_{\xi(\mathbf{x}) \geq T} P_{\mathbf{X}}(d\mathbf{x}).$$

In the Bayesian setting, the Bayes estimate of an unknown parameter in the sense of minimal mean squared error is simply the mean of the posterior distribution. Therefore, the value of $\mathbb{E}[P \mid \xi^n = \mathbf{g}^n]$ is an estimation of p . In other words, an estimator p_n of p is

$$p_n = \mathbb{E}[P \mid \xi^n] = \operatorname{argmin}_{Z \in \Delta} \mathbb{E}[(Z - P)^2], \quad (2)$$

where Δ is the set of square-integrable measurable function of ξ^n .

Note that the random variable P conditioned on ξ^n has the same distribution as the random variable P_n defined by:

$$P_n = \mathbb{P}(\xi_n(\mathbf{X}) \geq T \mid \xi_n) = \int_{\mathbb{X}} \mathbb{1}_{\xi_n(\mathbf{x}) \geq T} P_{\mathbf{X}}(d\mathbf{x}), \quad (3)$$

where ξ_n is the process ξ conditionally to $\xi^n = \mathbf{g}^n$.

For all $\mathbf{x} \in \mathbb{X}$, let $m_n(\mathbf{x})$ and $\sigma_n(\mathbf{x})$ be the mean and the variance of the random variable $\xi_n(\mathbf{x})$. We introduce the function

$$s_n(\mathbf{x}) = \mathbb{P}(\xi_n(\mathbf{x}) \geq T) = \Phi\left(\frac{m_n(\mathbf{x}) - T}{\sigma_n(\mathbf{x})}\right),$$

where $\Phi : \mathbb{R} \rightarrow [0, 1]$ is the cumulative distribution function of the standard Gaussian distribution. Thus, we have

$$p_n = \mathbb{E}[P_n] = \int_{\mathbb{X}} \mathbb{P}(\xi_n(\mathbf{x}) \geq T) P_{\mathbf{X}}(d\mathbf{x}) = \mathbb{E}[s_n(\mathbf{X})].$$

Finally, an estimator of p is

$$\hat{p}_n = \frac{1}{N} \sum_{i=1}^N s_n(\tilde{\mathbf{X}}_i),$$

where $\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_N$ is an N -sample of i.i.d. random variables with distribution $P_{\mathbf{X}}$.

The posterior mean value of P is easy to estimate and it does not require that one chooses an a priori distribution. However, its posterior distribution is untraceable. To learn about this distribution, a natural idea is to search its quantiles, i.e. to find for all $\alpha \in (0, 1)$ the number $F_{P_n}^{-1}(\alpha) \in [0, 1]$ such that

$$F_{P_n}^{-1}(\alpha) = \inf\{x \in \mathbb{R} : F_{P_n}(x) \geq \alpha\},$$

where $F_{P_n}(x) = \mathbb{P}(P_n \leq x)$. To achieve this, one can provide a quantile estimate build from realizations of the random variable P_n . However, this idea suggests for each simulation to get an exact realization of the process ξ_n and make an integration with respect to $P_{\mathbf{X}}$, which is very time consuming and could lead to numerical issues. See e.g [2].

Therefore, we propose to use an alternative random variable R_n in order to learn about the distribution of P_n .

4. An alternative estimation of the probability of failure

Let $U : \Omega \rightarrow [0, 1]$ be a random variable. Alternatively to P_n , we introduce the random variable R_n defined by

$$R_n = \mathbb{P}(s_n(\mathbf{X}) > U | U) = \int_{\mathbb{X}} \mathbb{1}_{s_n(\mathbf{x}) > U} P_{\mathbf{X}}(d\mathbf{x}). \quad (4)$$

Here $R_n = R_n(U)$ is a measurable function of U . For simplicity, we will write R_n instead of $R_n(U)$.

The distribution of R_n is one-dimensional and measurably depends on the distribution of U . It only involves the marginal laws of the process ξ_n , instead of the joint law for P_n . This random variable is easy to simulate because it does not involve realizations of the process ξ_n .

By introducing the alternative random variable R_n , our objective is to learn about the distribution of P_n .

Proposition 4.1. *Let us consider the random variables P_n and R_n defined in (3) and (4). The first moments of P_n and R_n are always equal if and only if the random variable U is uniform on $[0, 1]$.*

Therefore, by taking U uniform on $[0, 1]$, the mean value of R_n also provides an estimator of the probability of failure p . Now, let us recall the definition of the convex order between two random variables (see e.g. [7], Chapter 3).

Definition 4.1. The random variable X is said to be smaller than Y in the convex order, denoted $X \leq_{cx} Y$, if for all convex function φ

$$\mathbb{E}[\varphi(X)] \leq \mathbb{E}[\varphi(Y)],$$

provided these expectations exist.

Proposition 4.2. *Let us consider the random variables P_n and R_n defined in (3) and (4). The random variable P_n is smaller than R_n in the convex order if and only if U is uniform on $[0, 1]$.*

This result is useful to learn about the distribution of P_n . As mentioned in [7], the convex order leads to the following properties:

Proposition 4.3. *Let us consider the random variables P_n and R_n defined in (3) and (4). If $P_n \leq_{cx} R_n$, then*

$$(i) \text{ Var } [P_n] \leq \text{Var } [R_n],$$

$$(ii) \int_{\alpha}^1 F_{P_n}^{-1}(t) dt \leq \int_{\alpha}^1 F_{R_n}^{-1}(t) dt \text{ for all } \alpha \in [0, 1],$$

where $F_{P_n}^{-1}$ and $F_{R_n}^{-1}$ respectively denotes the quantile functions of P_n and R_n .

One can verify that the variance of P_n satisfies

$$\text{Var}[P_n] = \int_{\mathbf{x}} \int_{\mathbf{x}} \mathbb{P}(\xi(\mathbf{x}) \geq T, \xi(\tilde{\mathbf{x}}) \geq T) P_{\mathbf{x}}(d\mathbf{x}) P_{\mathbf{x}}(d\tilde{\mathbf{x}}).$$

The computation time requested by a Monte Carlo integration to estimate this variance can be high. Instead of this, we can use the variance of R_n to control the variance of P_n . Note that the variance of R_n is easy to estimate and only requires simulations of a uniform variable on $[0, 1]$ and an integration with respect to $P_{\mathbf{x}}$. Moreover, we show that

$$F_{P_n}^{-1}(\alpha) \leq \frac{1}{1-\alpha} \int_{\alpha}^1 F_{P_n}^{-1}(t) dt \leq \frac{1}{1-\alpha} \int_{\alpha}^1 F_{R_n}^{-1}(t) dt \leq \frac{\mathbb{E}[P_n]}{1-\alpha} = \frac{\mathbb{E}[R_n]}{1-\alpha}.$$

The upper-bound $\mathbb{E}[P_n]/(1-\alpha)$ is proposed in [2]. In the insurance and actuarial literature, the quantity $1/(1-\alpha) \int_{\alpha}^1 F_{R_n}^{-1}(t) dt$ is called the Conditional Value at Risk of R_n at level $\alpha \in (0, 1)$ and denoted by $\text{CVaR}_{\alpha}(R_n)$. In [1] and [4], the authors show that the Conditional Value-at-Risk CVaR_{α} of a random variable Z with a continuous distribution function F_Z is equal to the conditional expectation of Z given that $Z \geq F_Z^{-1}(\alpha)$, i.e.

$$\text{CVaR}_{\alpha}(Z) = \mathbb{E}[Z|Z \geq F_Z^{-1}(\alpha)]. \quad (5)$$

In fact, (5) is the usual definition of CVaR_{α} . Thus, the quantity

$$\text{CVaR}_{\alpha}(R_n) = \frac{1}{1-\alpha} \int_{\alpha}^1 F_{R_n}^{-1}(t) dt = \mathbb{E}[R_n|R_n \geq F_{R_n}^{-1}(\alpha)],$$

is a smaller upper-bound of the α -quantile of P_n than $\mathbb{E}[P_n]/(1-\alpha)$. It can easily be estimated, as long as α does not take values too close to 1.

For instance, let us consider the one-dimensional reliability case study presented in [3]. The real failure probability is $p = 0.2227$ and is computed thanks to a Monte Carlo integration with around 10^7 of simulations. The Gaussian process regression is performed with the R package **DiceKriging**. The estimation of mean value of the random variable P_n is 0.2274. The estimation of the mean value of R_n is 0.2275. By performing 200 realizations of the Gaussian process ξ_n , we provide estimations of the α -quantiles of P_n . All estimations are represented by the green line in The Figure 3, where α varies between 0 and 1. The blue line is the function $\alpha \mapsto \mathbb{E}[P_n]/(1-\alpha)$ and the red line is the function $\alpha \mapsto \text{CVaR}_{\alpha}(R_n)$. This figure shows that $\text{CVaR}_{\alpha}(R_n)$ is a better upper-bound of $F_{P_n}^{-1}(\alpha)$ than the one provide by Markov inequality.

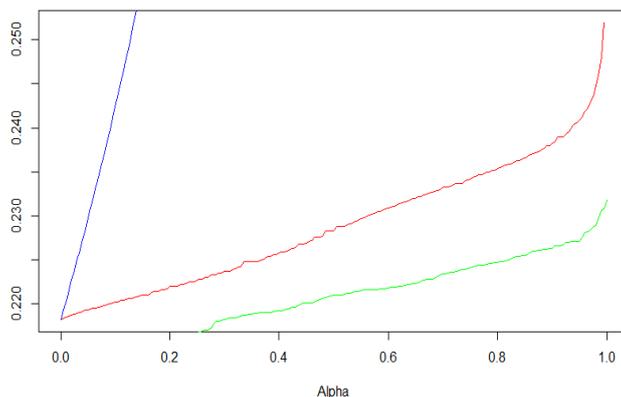


Figure 3. The blue line is the function $\alpha \mapsto \mathbb{E}[P_n]/(1-\alpha)$. The red line is the function $\alpha \mapsto \text{CVaR}_\alpha(R_n)$. The green line is the function $\alpha \mapsto F_{P_n}^{-1}(\alpha)$.

References

- [1] C. Acerbi and D. Tasche. On the coherence of expected shortfall. *Journal of Banking and Finance*, 26:1487-1503, 2002.
- [2] Y. Auffray, P. Barbillon, and J.-M. Marin. Bounding rare event probabilities in computer experiments. *Comput. Statist. Data Anal.*, 80:153-166, 2014.
- [3] J. Bect, D. Ginsbourger, L. Li, V. Picheny, and E. Vazquez. Sequential design of computer experiments for the estimation of a probability of failure. *Stat. Comput.*, 22(3):773-793, 2012.
- [4] V. Brazauskas, B. L. Jones, M. L. Puri, and R. c. Zitikis. Estimating conditional tail expectation with actuarial applications in view. *J. Statist. Plann. Inference*, 138(11):3590-3604, 2008.
- [5] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2006.
- [6] J. Sacks, T. J. Mitchell, W. J. Welch, and H. P. Wynn. Design and analysis of computer experiments. *Statist. Sci.*, 4(4):409-435, 1989. With comments and a rejoinder by the authors.
- [7] M. Shaked and J. Shanthikumar. *Stochastic Orders*. Springer Series in Statistics. Springer, New York, 2007.
- [8] M. Stein. *Interpolation of spatial data*. Springer Series in Statistics. Springer-Verlag, New York, 1999. Some theory for Kriging.