Computing conservative probabilities of rare events with surrogates

Nicolas Bousquet^{*}

Abstract

This article provides a critical review of the main methods used to produce conservative estimators of probabilities of rare events, or critical failures, for reliability and certification studies in the broadest sense. These probabilities must theoretically be calculated from simulations of (certified) numerical models, but which typically suffer from prohibitive computational costs. This occurs frequently, for instance, for complex and critical industrial systems. We focus therefore in adapting the common use of surrogates to replace these numerical models, the aim being to offer a high level of confidence in the results. We suggest avenues of research to improve the guarantees currently reachable.

Keywords: rare events, probability of failure, surrogates, accelerated Monte Carlo, probability upper bound, reliability, safety, guarantee, kriging, polynomial chaos expansion, neural network

1 Introduction

The estimation of probabilities of a rare event (often described as a critical failure) defined by

$$p := P(g(X) < y)$$

where $g: \Omega = [0, 1]^d \to \mathbb{R}$ is a deterministic function, $y \in \mathbb{R}$ such that p > 0 and X follows a uniform distribution on Ω , is the basic principle of many structural reliability studies [57]. Actually, X can be given more complicated distributions, but this framework remains very general. Difficulties occur when p can be very low and g can only be known through knowledge of a small number of realisations of Y = g(X), either by simulation of X, or from an observed sample of realizations of X collected in a possibly sequential manner. Simulation-based approaches are given particular consideration in what follows, as they provide a theoretical underpinning to purely empirical approaches. In this framework, many methods have been proposed to produce statistical estimators of p based on nqueries $x \mapsto g(x)$, with reduced asymptotic variance compared to the standard Monte Carlo (MC) approach, providing the well-known estimator p_n such that

$$\sqrt{n} (p_n - p) \xrightarrow{n \to \infty} \mathcal{N} (0, \sigma^2).$$

where $\sigma^2 = p(1-p)$ See [64] for a recent survey of such techniques. Their best representatives are based on importance sampling (IS) [23], possibly fed by tools from large deviation theory in high-dimensional settings [87, 83], sequential IS [73] and other adaptive IS methods as the cross-entropy method [35], directional sampling [65], line sampling [61, 34, 74], multilevel Monte Carlo [40, 50], variational approaches [90, 44] or importance splitting [55] (e.g., subset simulation [3, 4] and adaptive multilevel splitting [25, 20, 26]). If the best asymptotic variance that one may expect from these latter approaches is $p^2 \log(1/p) \ll \sigma^2$ [49], they require sampling multiple Markov chains, which can dramatically increase the number of queries $x \mapsto g(x)$ [14]. For this reason, many other approaches

^{*}SINCLAIR AI Laboratory, EDF R&D & LPSM, Sorbonne Université

are taking the gamble of building statistical surrogates (or meta-models) $x \mapsto \hat{g}_m(x)$ from the characteristics of $x \mapsto g(x)$ and a design of experiments (DOE) $\mathbf{x_m} = x_1, \ldots, x_m \in \Omega^m$, to reduce the overall cost of the calculation. Methods described in next paragraphs are summarized in Table 1.

1.1 Surrogate-based methods

A first class of methods considers a nonintrusive statistical approximation of $x \mapsto g(x)$ through quadratic response surfaces [45] or, much more commonly, in a Bayesian framework, through almost surely continuous Gaussian processes (kriging-based regression) whose mean defines the deterministic surrogate \hat{g}_m [81]. These imply assumptions of regularity (especially on the correlation structure [78]). Intrusive (adjoint-based) techniques can lead to define surrogates as *reduced order models* $(ROM)^1$, ie. defined using a reduced order basis [77], that may come with error bounds with respect to $x \mapsto g(x)$, as considered in [27].

Coming back to kriging-based surrogates, various sequential simulation techniques can then be used to benefit from the meta-model uncertainty, for example based on targeted Integrated Mean Square Error (tIMSE) [75] or Bayesian Stepwise Uncertainty Reduction (SUR) strategies [8], possibly combined with subset sampling methods [9]. Such strategies are defended by strong theoretical results [7]. Over the years, many adaptive variants of kriging-based methods have been proposed to estimate p (e.g. [15]). Active learning strategies, mixing sampling within Ω (MC, other better space filling techniques as Latin Hypercube Sampling, importance sampling, subset sampling, etc.) and kriging to produce estimates of p, have been particularly popularized. Let us cite the so-called AK-MCS methods (see [56] and [66] for recent reviews, and references therein). In engineering, kriging-based approaches appear useful to determine the subset of Ω leading to failures (e.g., [6, 63]), a dual gain in estimating p. Most famous competitors to Gaussian process meta-modeling are polynomial chaos expansions [82], implied within similar combinations of techniques to reach the estimation of very low probabilities. Both approaches have numerous merits in the general field of uncertainty quantification (e.g., versatility, easiness for conducting sensitivity analyses) and are the subject of major research aimed at reducing their computational cost in storage and inference (typ. $O(m^2)$ and $O(m^3)$ for kriging) and handling large dimensions d. Finally, the nesting of transformations in the regression function, for example through deep Gaussian processes [33], makes it possible to capture the complexities of the topological manifold defined by $x \mapsto q(x)$, at the cost of a large sample size. Neural networks and their universal approximation capability provide a final subclass of surrogates used in the latter context, aiming at diminishing the regularity assumptions. See [56] and [54] for two recent reviews of the various surrogate techniques used in the field of reliability, to complete this brief summary.

A second class of methods is concerned solely with the construction of a surrogate $\hat{\Gamma}_m$ of the failure (or *limit state*) surface

$$\Gamma = \{x \in \Omega, g(x) = y\}$$

which is defined under continuity assumptions, and considering that the problem is a perfectly separable binary classification problem. This generalizes in finding excursion sets, as described hereinafter. Engineering methods like FORM/SORM [37] make the assumption that Γ can be locally approximated by linear or quadratic surfaces, and transform the estimation problem into an optimization problem. They are usually combined with IS, as techniques based on large deviation theory [87]. Combined with subset simulation, Support Vector Machines (SVM) or combinations of SVM were proposed by [17] and [19] to approximate Γ (see [80] for a review), while neural networks were preferred by [72, 18, 60]. Kriging-based classification was proposed by [38], while [59] chose polynomial chaos expansion, accompanied with cross-entropy importance sampling [58].

¹As recalled by [27], such methods are known to be "very efficient for the numerical approximation of problems involving the repeated solution of parametric Partial Differential Equations (PDEs)".

These approximations are often combined, again, with sequential sampling strategies. The AK-MCS methods evoked hereinbefore fall into this category too, as they are underlyingly based on a classifier (approximating how a vector x is far from Γ) defined by a kriging predictor.

1.2 Ensuring conservatism

When it comes to the real-life estimation of p to characterize the safety of critical systems (e.g., nuclear/aeronautical/spatial structures and processes), where g is typically a certified numerical model, the question arises of the credit to be given to estimators based on \hat{g}_m or $\hat{\Gamma}_m$. Do they provide "good", "sufficient" estimators of the true probability p? Can we establish insurance rules based on these estimators? Despite the continuous improvement of these methods and their combinations, and the subsequent reduction of the error attributable to the use of surrogates² – in particular by adaptive sampling techniques and selected queries $x \mapsto g(x)$ [5, 48, 58] – most surrogate-based methods are still lacking of theoretical guarantees adapted to the context. The asymptotic convergence to p of any estimator of the kind

$$\hat{p}_{n,m} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{\hat{g}_m(x_i) < y\}} \qquad (MC \text{ or } IS \text{ with surrogate of } g) \tag{1}$$

or

$$\hat{p}_{n,m} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{x_i < \hat{\Gamma}_m\}}, \qquad (MC \text{ or } IS \text{ with surrogate of } \Gamma)$$
(2)

where < is some partial ordering rule (allowing binary classification), appears *weak*, in the sense it does not answer to the requirements that can be expected from certification authorities³. While the consistency of (1) and (2) in (m, n) is of course necessary, this asymptotic property remains a figment of the imagination when g is expensive, and/or d is large. And, more generally, because we can hardly define the reality of an asymptotic regime, even with a central limit theorem. Stronger guarantees should be non-asymptotic, for instance provided through concentration inequalities (or PAC-Bayes-type inequalities) as

$$P\left(\left|\hat{p}_{n,m} - p\right| > \varepsilon\right) \leqslant \alpha_{m,n,\varepsilon} \tag{3}$$

where $\alpha_{m,n,\varepsilon}$ remains ideally very low for achievable (m,n) (at the order of p) and small ε and $\alpha_{m,n,\varepsilon} \to 0$ when $(m,n) \to \infty$ for any ε .

Other strong guarantees are likely to come from the specific choice of ROM, as they often come with error controls with respect to g. Evoked in the recent work [27], such properties seem not to have been yet too much explored, maybe because such ROM are typically made to represent regular behaviors rather than extreme ones. Finally, other strong guarantees are related to robustness to variations on g (or, similarly, on the real distribution of X, through a transport from the uniform distribution introduced hereinbefore). Even most of MC-based methods provide estimates of p that lack of robustness, as their coefficient of variation goes to infinity when p goes to 0 [46]. This type of problem involves both robustness studies and uncertainty quantification. Robust analysis theories struggling against the probabilistic misspecification of epistemic uncertainties, such as the info-gap theory [10, 11], which uses convex epistemic models⁴, are used in *hybrid* reliability studies [92].

$$X \sim \underbrace{\text{distribution } P}_{uncertain} \xrightarrow{Consequence} p \in (p_n^-, p_n^+)$$

 $^{^{2}}$ A part of this error being linked to dimension-reduction, see [54].

³As said in [39], "Certification authorities will most likely require safeguards", while the authors of [54] insist: "The expectation is inappropriate for risk-averse stakeholders, as it does not capture any notion of variability and does not quantify rare events"

 $^{^{4}}$ But can be associated to other epistemic models (e.g. possibility distributions, p-boxes, Dempster-Schafer structures...).

Using tools derived from random set theory [88], they can be used to define bounds on probabilities of failure. See [1] for more information. However, computing these bounds is itself subject to the same computational problems as that of a single probability p for which the model $X \sim \mathbb{P}$ is well defined, and usually requires reduction variance techniques or/and surrogates. An alternative is to make geometric assumptions on the distribution of the output g(X) to get conservative bounds, in the spirit of [85]. These authors provide rules in the specific case of survival analysis, where prior information on lifetime distributions is available.

One way of simultaneously addressing the previous concerns is to produce *conservative* estimators of p, ie. upper bounds $p_{n,m}^+ \ge p$, such that $p_{n,m}^+ \to p$ and $p_{n,m}^+ - p$ be not too wide [18]. Ideally, one would like the order of magnitude of $p_{n,m}^+$ to be similar to that of p: if $p \sim 10^{-q}$, then it is hoped to get $10^{-q} \le p_{n,m}^+ \ll 10^{-q+1}$. Note that if some algorithm can be run to determine $p_{n,m}^+ \in [p, 1)$, by symmetry the same algorithm can be adapted to determine a lower bound $p_{n,m}^- \in (0, p]$, provided p > 0.

Noticing that $p = \operatorname{Vol}(\Omega_y)$ where Ω_y is the so-called *excursion set* (or *level set*) [16, 6]

$$\Omega_y = \{ x \in \Omega : g(x) < y \},\$$

any upper bound $p_{n,m}^+$ for p corresponds to the volume of a subset $\Omega_{y,m,n}^+$ such that $\Omega_y \subseteq \Omega_{y,m,n}^+ \subsetneq \Omega$. The complementary subset $\overline{\Omega}_{y,m,n}^+ = \Omega/\Omega_{y,m,n}^+$ corresponds to a set smaller than the safe set $\overline{\Omega}_y = \{x \in \Omega : g(x) \ge y\}$. When the elements of $\overline{\Omega}_{y,m,n}^+$ are included in $\overline{\Omega}_y$ with a large probability β , such sets $\overline{\Omega}_{y,m,n}^+$ are called *conservative* by [43, 16]. Still in a Bayesian setting, these latter authors proposed a kriging-based approach to estimate this conservative set (thanks to the posterior distribution of g) from a fixed DOE. This work was recently extended by [6], who proposed adaptive strategies based on specific SUR criteria to reduce the uncertainty of this estimator. While offering excellent performance in terms of volume recovery on low-dimensional cases, this approach costs ~ 200 queries $x \mapsto g(x)$ on a real industrial example in 2D. While it has not been used specifically to compute upper bounds, it partially answers to the questions above, as it is conditioned by the relevance of the choice of the Gaussian process.

To the best of our knowledge, obtaining appropriate guarantees for the certification of a calculation of a probability of failure using surrogates therefore remains a tough problem. The aim of this article is to summarize some useful results for further research in this area. First (Section 2), we recall some strong conservatism results linked to the use of certain geometric properties of g. These are difficult to check by nature, but can be much more easily used to select certain classes of surrogates. They provide theoretically deterministic bounds on a probability p calculated from a surrogate. We also consider in Section 3 approaches inspired by [39] and [54], which seek to "bias" a surrogate to ensure a form of conservatism in the estimation of p. A discussion of future prospects concludes this work, focusing on the obstacles to be overcome.

2 Obtaining deterministic bounds on p from geometrical properties

The strongest desirable conservatism condition consists in producing deterministic upper bounds p_n^+ on p, independently on any additional m-sample that could be used to estimate a surrogate. The only known bounds are intrinsically linked to geometrical conditions, sometimes induced by conditions of regularity which we recall and discuss below. For the sake of generality we denote by \mathbb{P} the probability measure associated to X over a subset of Ω (which, in our context, is simply the uniform measure).

Table 1: A summary	of most popular	methods	dedicated	to the	e estimation	of rare event	
probabilities.							

Variance reduction methods		(survey: [64])
Importance sampling (IS)		[23]
	+ large deviation theory	[87, 83]
	+ spectral decomposition	[95]
	+ variational approaches	[90, 44]
	Sequential IS	[73]
	Cross-entropy IS	[35]
Directional sampling		[65]
Line sampling		[61, 34, 74]
Multilevel Monte Carlo		[40, 50]
Importance splitting		[55]
	Subset simulation	[3, 4]
	Adaptive multilevel splitting	[49, 25, 20, 26]
Methods based on a surrogate of g		(survey: [56])
Reduced order models		[77]
Quadratic response surfaces		[77] [45]
Mean predictors of Gaussian processes		[±0] [81]
Mean predictors of Gaussian processes	\pm Sequential strategies	[01] [7]
	Ex 1: Bichon criterion	[1] [15]
	Ex 1: iMSE criterion	[75]
	Ex 2: Stepwise Uncertainty Reduction (SUR)	[8]
	Ex 3 : SUB + Subset sampling	[9]
	Ex 4 : Active learning (AK-MCS)	[56, 66]
	Ex 5 : derived from inversion methods	[63]
	Ex 6 : derived from excursion sets	[16, 6]
Polynomial chaos expansions		[82]
		LJ
Methods based on a surrogate of Γ		(survey: [80])
		[]
Linear response surfaces (FORM)		[37]
Quadratic response surfaces (SORM)		[37]
	+ 1S	[21]
Support Vector Machines (SVM)		[17]
	Combination of SVM	[19]
Neural networks		[72, 18]
	Adaptive strategies	[U0] [20]
Adaptive kriging		[38] [50]
Polynomial chaos expansion	areas entrony IS	[59] [59]
	\pm cross-encropy is	ျပ၀

2.1 Lipschitz smoothness

Theorem 1 is derived from a sequence of results recently obtained by [14]. It uses classical tools in harmonic analysis: dyadic cubes, for which a constructive definition must first be provided, illustrated by Figure 1.

Definition 1 (Dyadic cubes.). Dyadic cubes in $\Omega = [0,1]^d$ are a (possibly infinite) collection of cubes

$$\begin{cases} Q_0, \\ Q_{1,1}, \dots, Q_{1,2^d}, \\ Q_{2,1,1}, \dots, Q_{2,1,2^d}, Q_{2,1,1}, \dots, Q_{2,2,2^d}, \dots, \dots, Q_{2,2^d,1}, \dots, Q_{2,2^d,2^d} \\ \dots \end{cases}$$

such that:

- $Q_{j,\ldots}$ has sidelength 2^{-j} for $j \in \mathbb{N}$, with $Q_0 = [0,1]^d$;
- for any $j \in \mathbb{N}$, the $Q_{j,\dots}$ define a partition of Ω ;
- each $Q_{j,...}$ has 2^d children cubes $Q_{j+1,...}$ build by performing a 2^d split of $Q_{j,...}$;
- each $Q_{j,...}$ has exactly one parent cube, for j > 0.

A cube Q with sidelength 2^{-j} is said to have a depth j.





Each cube Q is then determined by its depth j and its center, noted c_Q . Reusing the notations of [14], it can be labelled by a query to $x \mapsto g(x)$ as follows:

- $Q \in \mathcal{I}$ if $g(c_Q) > y + L2^{-j-1}$ (inside the safe subset of Ω),
- $Q \in \mathcal{O}$ if $g(c_Q) < y L2^{-j-1}$ (outside the safe subset of Ω),
- $Q \in \mathcal{U}$ otherwise.

Given a choice of maximal depth k, a recursive algorithm given in [14] allows n cubes to be labelled, as precised in the following theorem. Denote then \mathcal{I}_n the set of all cubes labelled as \mathcal{I} along this algorithm (and similarly let us denote \mathcal{O}_n and \mathcal{U}_n).

Theorem 1 (Derived from [14].). Under the following assumptions:

- (i) The distribution of X on Ω admits a bounded density function with respect to the Lebesgue measure λ;
- (ii) The function g is assumed to be L-Lipschitz with respect to the supremum norm on \mathbb{R}^d , ie.

$$|g(x) - g(\tilde{x})| \leq L ||x - \tilde{x}||_{\infty}, \quad x, \tilde{x} \in \Omega;$$
(4)

(iii) There exists a constant C > 0 such that

$$\lambda\left(\left\{x\in\Omega: |g(x)-y|\leqslant\delta\right\}\right) \leqslant \frac{C}{L}\delta, \quad \delta>0 \quad (level \ set \ condition); \tag{5}$$

then

$$p_n^- \leqslant p \leqslant p_n^+ \tag{6}$$

where

$$p_n^+ = 1 - \sum_{Q \in \mathcal{I}_n} \mathbb{P}(X \in Q) \quad and \quad p_n^- = p_n^+ - \sum_{Q \in \mathcal{U}_n} \mathbb{P}(X \in Q),$$

with, for $d \ge 2$,

$$n \leq 4C2^{d-1}$$

and $p_n^+ - p_n^- \leq Cn^{-\frac{1}{d-1}}$,

this latter result being the optimal convergence rate, in the sense it cannot be improved by any other algorithm defined under the sole general assumptions.

In our context, Assumption (i) is always true, and note that $\mathbb{P}(X \in Q) = 2^{-dj}$ for any Q of sidelength j. Notice that [14] consider more broadly that X can follow any distribution on Ω provided (i) remains true. In this case, the computation of the $\mathbb{P}(X \in Q)$ is not trivial and requires numerical methods, that bring stochastic errors on the estimation of bounds (p_n^-, p_n^+) . Using a splitting approach, the authors can still control the bounding of p and the convergence of the difference of estimated bounds towards 0. To our knowledge, these results were only applied to a univariate toy case with failure probability of the order 10^{-3} (see Table 4).

As explained by the authors, the level set Assumption (iii) reflects the fact that g is not too much flat in the vicinity of the limit state (or level set, or failure surface) $\Gamma = g^{-1}(y)$. This condition appears mild under a continuous differentiability assumption on g. Besides, the knowledge of M is not necessary for the algorithmic task.

Estimating (or bounding) the Lipschtiz constant L is however required and is a more serious difficulty for carrying out the real computation of the bounds. This is generally a difficult problem, as it involves being able to calculate a maximum distance between the outputs of g. However, if g is a surrogate chosen as a neural network, owning the Lipschtiz property, the tight estimation of L has become over the recent years a challenge for which constructive solutions have recently been published. Especially, the approach proposed by [41], that estimates L by solving a semidefinite problem, has become very popular in the machine learning community. See besides [71] for the estimation of Lipschitz constants of monotone deep equilibrium models.

2.2 Monotonicity

A particular case of functions $x \mapsto g(x)$ are the monotonic ones, when monotonicity is understood with respect to the following partial order on \mathbb{R}^d defining Pareto dominance :

Definition 2. Let $x = (x_1, \ldots, x_d)$ and $x' = (x'_1, \ldots, x'_d) \in \mathbb{R}^d$. If $x_i \leq x'_i$ for all $i = 1, \ldots, d$, x is said to be dominated by x', and it is denoted $x \leq_P x'$.

To simplify things in the following, the monotonicity of g is assumed to be similar for all dimensions in Ω : g is supposed to be globally increasing (possibly at the price of a reparameterization), ie; $\forall \{x, x'\} \in [0, 1]^d \times [0, 1]^d$ such that $x \leq_P y'$, then $g(x) \leq g(x')$. An immediate consequence is that each $x \leq_P x'$, when $g(x') \leq y$, is such that $g(x) \leq y$. This led [36] and [18] then [68] and [19] to define a class of algorithms that simultaneously produce deterministic bounds (p_m^-, p_m^+) and consistent statistical estimators $\hat{p}_{n,m}$ of p, from the Lebesgue measures (volumes) of sequences of nested random dominated sets

$$\mathbb{U}^{-}(A) = \bigcup_{x \in A \cap \mathbb{U}^{-}} \{ u \in [0,1]^{d} : u \le x \},$$
$$\mathbb{U}^{+}(A) = \bigcup_{x \in A \cap \mathbb{U}^{+}} \{ u \in [0,1]^{d} : u \ge \mathbf{x} \}.$$

with

$$\mathbb{U}^- = \{ x \in [0,1]^d : g(x) \leq y \} \text{ and } \mathbb{U}^+ = \{ x \in [0,1]^d : g(x) > y \}.$$

Figure 2 offers an 2D-illustration of the fundamentals of this method. The bounds obtained this way were used to guarantee the result of safety analysis for some industrial case-studies (e.g., [86]). Sequential surrogates of the limit state surface Γ can be used to guide new design points to lower the volume $p_m^+ - p_m^-$ (and jointly improve $\hat{p}_{n,m}$), but the core of the approach stands on the gain yielded by the deterministic (upper) bounds, which can be computed exactly in low dimensions [18] or else by specific numerical methods [19].



Figure 2: An illustration in dimension 2 of the bounding principle due to monotonicity. On the left is a set A of m points included in $[0,1]^d$. On the right, we plot the dominated spaces $\mathbb{U}^-(A)$ and $\mathbb{U}^+(A)$ surrounding Γ . It is clear that $\mathbb{U}^-(A) \leq_P \Gamma \leq_P \mathbb{U}^+(A)$. The volume of these subspaces allows to compute deterministic bounds for p: $\lambda(\mathbb{U}^-(A)) \leq p \leq 1 - \lambda(\mathbb{U}^+(A))$, with λ the Lebesgue measure.

Actually, for building such bounds it is enough to that the limit state surface Γ be p-monotone in the following sense:

Definition 3. Denote μ the Lebesgue measure on \mathbb{R}^d . Let $A \subset [0,1]^d$. Define the dominated sets

$$\mathbb{V}^{-}(A) := \bigcup_{x \in A} \{ u \in [0,1]^d : u \leq_P x \},\$$
$$\mathbb{V}^{+}(A) := \bigcup_{x \in A} \{ u \in [0,1]^d : u \geq_P x \}.$$

Let $\alpha \in (0,1)$ and S be a set in $[0,1]^d$. The set S is said α -monotonic if for all $u, v \in S$ such that $u \neq v$, u is not strictly dominated by v and if $\mu(\mathbb{V}^-(S)) = \alpha$.

As a consequence of the monotonicity of g (Proposition 5.1 in [19]), a major feature of Γ is its *p*-monotonicity, provided that Γ is simply connected and $\lambda(\Gamma) = 0$. Note that these two last assumptions are required in [19] to estimate consistently Γ , interpreted as the frontier of a separable classification problem.

Monotonic sub-behaviors are often key to understand complex behaviors summed up by some simple rule as testing if g(X) < y, as demonstrated by Bshouty's monotone theory for learning Boolean functions [22], many works related to monotonicity detection in such functions (e.g. [47]) and the importance of monotonicity in causal analysis and SAT theory. Modeling and sensitivity analysis theories make extensive use of tools to detect monotonicity between inputs and X and outputs Y (typically via correlation coefficients). The exhibition of partial or total monotonicity properties of a model or phenomenon $x \mapsto g(x)$ defines a semantics allowing its interpretability and the obtaining of guarantees on the predictive outcomes g(x) [84]. Using monotonicity property from expert knowledge can help to produce more suitable surrogates in engineering prblems [31, 53, 32]. Hence monotonicity facilitates the use of g in critical domains such as clinical testing (e.g., [69]) or credit scoring (e.g., [76]). See the previous references for various industrial illustrations.

For these reasons, numerous models that can be used as surrogates of g or Γ present global or local monotonicity properties, from the simplest one (linear model) to complex ones (e.g., deep lattice networks [94] or GAMI-nets [93], that allow *monotonization*). See for instance [24] for a review of monotonic classification.

For instance, useful generalizations of linear models are the m-Continuous Piece-Wise Linear functions (m-CPWL) which can be interpreted as piecewise monotonic surrogates:

Definition 4. [29, 70, 28] A function $g: \Omega \to \mathcal{Y}$ is a m-Continuous Piece-Wise Linear function (m-CPWL) is there exists K finite sets of disjoint complex polytopes $A_{k_{k=1}}^m$ such that $\bigcup_{k=1}^m A_k = \Omega$ and g restricted to the domain A_k , denoted as $g_{|A_k}: A_k \ni x \mapsto g(x)$ is affine for each $k \in \{1, \ldots, m\}$.

This versatile class encompasses neural networks with piece-wise linear activations such as ReLu or hard tanh, that correspond to $\max(0, x)$ and $\max(-1, \min(1, x))$, respectively. As explained in [2] (Chapter 3), feedforward neural networks can be described as piece-wise linear functions that divide the input space into multiple linear, where the network itself behaves as an affine function within each region [52, 51, 28]. Reusing the notations in Theorem 2.2 from [2], we describe the hyperrectangles A_k by

$$A_k = \bigotimes_{i=1}^p A_{i,k}$$

where $A_{i,k} = [l_{i,k}, r_{i,k}]$ with $(l_{i,k}, r_{i,k}) \in \overline{\mathbb{R}}^2$. Each component $g_{|A_k|}$ is represented by

$$g_{|A_k}(x) = \sum_{i=1}^p a_{i,k} x_i + b_k$$

where the $(a_{i,k}, b_k)$ are real numbers. Consequently,

$$g(x) = \sum_{k=1}^{m} \left(\sum_{i=1}^{p} a_{i,k} x_i + b_k \right) \mathbb{1}_{A_k}(x)$$
(7)

which is clearly piecewise monotonic.

Numerical experiments conducted with the following toy example proposed in [18] were conducted until dimension 15 in [13]. Approximate nested uniform sampling within the non-dominated space (defined by $\Omega/(\mathbb{U}^+(A) \cup \mathbb{U}^-(A))$) in Figure 2) was produced using a semi-adaptive MCMC summarized in Appendix A.

Example 1. Given a fixed dimension d, denote

$$V_d = \tilde{g}_d(Z) = Z_1 / (X_1 + \sum_{i=2}^d Z_i)$$

where the Z_i are Gamma distributed: $Z_i \mathcal{G}(i+1,1)$, for i = 1, ..., d. The function \tilde{g}_d is increasing in Z_1 and decreasing in all the variables Z_i for $i \ge 2$. The same monotonicity can be deduced for $x \mapsto \tilde{g}_d \circ T^{-1}(x)$ where $x \in \Omega$ and T is obtained from the cumulative distribution functions of Z_i : $T = (F_{Z_1}, ..., F_{Z_d})$. The variable V_d follows the Beta distribution $\mathcal{B}(2, 2^{-1}(d+1)(d+2) - 3)$. For $p \in [0, 1]$, let $q_{d,p}^1$ be the quantile of order pof Z_d , and let $g_d(x) = \tilde{g}_d \circ T^{-1}(x) - q_{d,p}$. Hence, with y = 0, $P(g_d(X) < 0) = p$ for all $d \ge 2$.

We can then compare the theoretical value p with bounds (p_n^-, p_n^+) obtained through the nested algorithmic approach evoked above, after n = 200 steps, and a basic one based on "brute force Monte Carlo". Relative precisions $(p_n^+ - p_n^-)/p$ are provided on Table 2 for several dimensions d, while Figure 3 displays how the computing time varies in function of d, for $p = 5.10^{-4}$.

Table 2: Mean values of relative precision $(p_n^+ - p_n^-)/p$ after 200 iterations for the MCMC-based method.

		Brute force MC	MCMC
d	p		
2	5.10^{-2} 5.10^{-3} 5.10^{-4}	$\begin{array}{c} 0.233 \\ 0.304 \\ 0.06 \end{array}$	$0.208 \\ 0.276 \\ 0.04$
3	5.10^{-2} 5.10^{-3} 5.10^{-4}	$1.18 \\ 1.78 \\ 2.59$	$1.09 \\ 1.66 \\ 2.65$
4	5.10^{-2} 5.10^{-3} 5.10^{-4}	$2.53 \\ 5.57 \\ 8.98$	$2.38 \\ 5.40 \\ 8.59$

Obviously, the computational time is significantly decreased using a MCMC approach, for the same order of magnitude for the highest dimensions. However, reaching small relative orders of magnitude remains difficult even for dimension 4, with $p_n^+ \sim 5p$, which fits with the results previously obtained by [68]. This remains a bit frustrating for the industrial practice, notably, but not hopeless. More simulations, and more clever simulations, for instance guided by a surrogate of Γ sequentially updated, as the one proposed in [19], should be required to reach $p_n^+ \sim p$.

Convexity and quasi-convexity. It must be noticed that convexity and quasi-convexity properties on g or Γ can certainly help achieve this ambition. To our knowledge, the use of such properties has still little explored, apart in [67], and in the literature dedicated to computing probabilities related to PDE solving ; see for instance [89] for a summary presentation. Available results highlight speed orders on the distance between Γ and surrogates of Γ , as exponential functions of the dimension; this distance can be bounded on convexity arguments [89]. However, bounding the implied relative error of the probability of failure still largely remains an open problem.



Figure 3: Average calculation time (in seconds) as a function of size. "No MCMC" means "brute force Monte Carlo." Figure extracted from [13].

3 Biasing surrogates to get conservative failure probabilities

3.1 Principle

When the use of a surrogate \hat{g}_m trained from a *m*-sample is required for computational reasons, the estimation of *p* by a $\hat{p}_{n,m}$ should be ideally such that

$$p \leqslant \hat{p}_{n,m} \tag{8}$$

with a large probability $1 - \alpha$, ideally with $\alpha = 0$. There are arguments to understand that α can possibly be strictly positive: (a) the computations are made using a fixed data set (part of which is used to train the surrogate); it is therefore conceivable that not all the uncertainty about $\hat{g}_m(X)$ can be taken into account; (b) when the dimension increases, as said in [54], "surrogate modeling techniques are often built on lower dimensional subspaces identified by dimension-reduction techniques". These techniques introduce an error that is difficult to quantity.

To diminish α and more generally enforcing the so-called *first-order stochastic domi*nance constraint [54]

$$P(\hat{g}_m(X) < y) \geq P(g(X) < y) \quad \forall y \in \mathbb{R},$$
(9)

an idea originated from [39] and also exploited by [54] is to "bias" (or "shift") the surrogate, replacing $x \mapsto \hat{g}_m(x)$ by

$$x \mapsto \hat{g}_m(x) + \theta$$

to ensure conservativeness. Obviously, in our context θ should be chosen such that (8) is verified with $\hat{p}_{n,m}$ having the closest possible order of magnitude than the one of p. Obtaining too much difference in terms of order of magnitude should lead to modify the choice of the surrogate.

3.2 Using Bernstein concentration inequalities

More precisely, [39] prove the following result, using a uniform Bernstein-type inequality:

Theorem 2 (Inspired from [39], Corollary 1). Denote \hat{g}_m a surrogate of g trained over a m-sample and consider an independent test set of $n \ge 2$ iid values X_1, \ldots, X_n drawn over Ω . Denote

$$\theta^* = \min\left(0, \min_{1 \le i \le n} \left\{ \hat{g}_m(X_i) - g(X_i) \right\} \right).$$

Then, with probability at least $1 - \alpha$ over the choice of the test set, there exists a strictly positive constant C < 6 such that

$$P(\hat{g}_m(X) + \theta^* > g(X)) \leqslant B(n, \alpha) := \frac{C}{n} \log(n/\alpha)$$

We can straightforwardly derive from it the following result, assuming $B(n, \alpha) = p$. Corollary 1. Under the assumptions of Theorem 2,

$$\mathbb{P}\left\{P\left(\hat{g}_m(X) + \theta^* \leqslant g(X)\right) \ge 1 - p\right\} \ge 1 - \lambda(n, p) \tag{10}$$

with

$$\lambda(n, p) = \min(1, n \exp(-np/C))$$

Then, with probability at least $1 - \lambda(n, p)$ over the choice of the test set,

$$\hat{p}_m := P(\hat{g}_m(X) + \theta^* < y) > p.$$

While it is interesting to note that this result does not depend on the dimension of X, a large number of values n is required to produce a nontrivial lower bound in (10), for low probabilities p. This is exemplified by plotting representatives values of $\lambda(n, p)$, setting C = 6, in Figure 4, and by computing the value of n such that $\lambda(n, p) \sim p$, displayed (on the log scale) on Figure 5. Typically, one needs $n \simeq 512$ to obtain $\lambda(n, p) = p = 10^{-1}$ then $n \simeq 8280$ to get $\lambda(n, p) = p = 10^{-2}$. In practice, this cost remains probably too high for industrial applications. The authors [39] have logically proposed a complete procedure for learning θ in parallel of the surrogate parameters, which is enforced in the following numerical experiments.

Reusing the toy example (1), numerical tests were conducted with neural networks surrogates. They were considered as good enough and retained for the study if both their accuracy and their Q^2 predictivity coefficient [42] were estimated above 90% and 0.9, respectively. Neural networks were simple feedforward networks with logistic activation functions, 2 to 3 neuronal layers, each layers having 2 to 4 neurons. On Table 3, some results of these numerical tests are summarized.

These results confirm the theoretical reservations set out hereinbefore; it seems unacceptable to obtain, for "reasonable" low probabilities, upper bounds that can almost vary by an order of magnitude when the dimension remains low, and which turn out to be false upper bounds in a significant number of cases.

3.3 Learning minimal bias

A close approach but specifically dedicated to obtaining upper bounds on probabilities was then proposed by [54], inspired by the work [91], who build risk-averse surrogate models using stochastic dominance. Given a training sample $(x_i, y_i)_{1 \le i \le m}$ and a parametric surrogate $x \mapsto g_{\eta}(x)$ mimicking $x \mapsto g(x)$, they pose the problem of learning (θ, η) by choosing a specific weighted loss/cost function

$$\min_{\theta,\eta} \sum_{i=1}^{m} \omega_i \left(y_i - g_\eta(x_i) - \theta \right)^2 \tag{11}$$

to minimize under the first-order stochastic dominance constraints

$$\begin{cases}
\sum_{i=1}^{m} \omega_{i} \mathbb{1}_{(-\infty,0]} \left(g_{\eta}(x_{i}) - g_{\eta}(x_{1}) \right) &\leq \sum_{i=1}^{m} \omega_{i} \mathbb{1}_{(-\infty,0]} \left(y_{i} - g_{\eta}(x_{1}) - \theta \right), \\
\dots &\leq \dots \\
\sum_{i=1}^{m} \omega_{i} \mathbb{1}_{(-\infty,0]} \left(g_{\eta}(x_{i}) - g_{\eta}(x_{m}) \right) &\leq \sum_{i=1}^{m} \omega_{i} \mathbb{1}_{(-\infty,0]} \left(y_{i} - g_{\eta}(x_{m}) - \theta \right)
\end{cases}$$
(12)



Figure 4: Values of $\lambda(n,p)$ with C = 6, in function of $\log(n)$, for some values of p.



Figure 5: Values of log(n) such that $\lambda(n,p) \simeq p$, for some typical values of p.

Table 3: Some applied results on the reality of the conservative assessment of p = P(g(X) < 0) using simple neural networks, for the toy example (1), by the shifting/biasing approach originally proposed by [39] for creating conservative surrogates. The two last columns present the average conservative estimate \hat{p}_n of p produced by the surrogate, and the probability that a particular estimate actually be a wrong upper bound for p. These results are produced using $\lfloor 100/(p(1-p)) \rfloor$ repetitions of the procedure, and training datasets of length 50d to ensure a comparable precision of results.

			Surrogate features		
p	d	Q^2 without imposed bias	Q^2 with imposed bias	\hat{p}_n	$\mathbb{P}(\hat{p}_n < p)$
10^{-1}	2	0.99	0.97	$1 \ 12 \ 10^{-1}$	0.08
10	3	0.98	0.94	$1.23.10^{-1}$	0.10
	4	0.95	0.92	$1.31.10^{-1}$	0.13
	5	0.94	0.92	$1.38.10^{-1}$	0.17
10^{-2}	2	0.98	0.97	$1.9.10^{-2}$	0.09
	3	0.96	0.94	$2.4.10^{-2}$	0.14
	4	0.94	0.91	$3.1.10^{-2}$	0.21
10^{-3}	2	0.96	0.94	$2.8.10^{-3}$	0.18
	3	0.95	0.91	$3.9.10^{-3}$	0.25
	4	0.95	0.90	$6.2.10^{-3}$	0.34

This setting is the empirical version (ie., based on a given finite training dataset) of the estimation problem

$$\min_{\substack{\theta,\eta \\ \text{subject to } \theta + g_{\eta}(X) \ge Y}} \left\{ \left\{ (Y - g_{\eta}(X) - \theta)^{2} \right\} \right\}$$
(13)

where $Y_1 \geq Y_2$ means that Y_1 dominates Y_2 with respect to the first stochastic order, ie. when

$$P(Y_1 \leq t) \leq P(Y_2 \leq t) \quad \forall t \in \mathbb{R}.$$

Because of the discontinuity introduced by the indicator functions in (12), the authors[54] consider a continuous relaxation (originated from [30]) of the constrained problem (11-12), rewritten besides as a mixed integer optimization problem. Despite technical difficulties related to the choices of relaxations, the authors prove the relevance of this approximation and the conservativeness of the overall approach, ie.

$$P(\hat{g}_m(X) + \theta^* < y) \ge P(g(X) < y) \quad \forall y \in \mathbb{R},$$

where $\hat{g}_m = g_{\eta^*}$, (η^*, θ^*) being the solution of (13) and its empirical approximations. They also prove its applicability for several uni- and multidimensional examples using polynomial chaos expansion, and show the feasibility of choosing θ dependent on x. In the same paper [54], the authors defend another approach, based on the *risk quadrangle* [79], for building surrogates that conservatively estimate a specific risk measure associated with the subjective preferences of a stakeholder.

4 Discussion

This rapid review of tools and methods designed to support the construction and use of conservative estimators of probabilities of failure can be finally illustrated by a some typical numerical results, displayed in Table 4.

Table 4: Summary of some typical orders of magnitude of upper bounds for p found in the dedicated literature ("un" means unpublished), considering toy models and more realistic examples. The value n is the typical MC number of queries (training samples) $x \mapsto g(x)$ required for the computation (or samples for training a surrogate). The number of surrogate runs for a MC analysis is not precised, as it can be chosen as large as wished. Results $\mathbb{P}(p > p_n^+)$ are averaged using validation samples (from 100 to 9000 values) if surrogates are used.

p	d	n	upper bound p_n^+	$\mathbb{P}(p > p_n^+)$	use case	details	references
10^{-1}	10	30	$1.02.10^{-1}$	0	wing weight	biased surrogate-based (1D active space)	[54]
5.10^{-2}	3	100	$5.2.10^{-2}$	0	toy model	biased surrogate-based	[54]
5.10^{-2}	10	100	10^{-1}	0	truss structure	biased surrogate-based	54
$2.1.10^{-3}$	1	32	$2.5.10^{-3}$	0	toy model	no surrogate,	[14]
						Lipschitz constant known	
10^{-3}	6	300	$1.2.10^{-2}$	0	toy model	no surrogate, g monotonic	[68]
10^{-4}	3	200	2.10^{-4}	0	toy model	no surrogate, g monotonic	[68]
10^{-4}	4	200	3.10^{-4}	0	hydraulics	no surrogate, g monotonic	un.
10^{-4}	5	250	6.10^{-4}	0	toy model	no surrogate, g monotonic	[68]
5.10^{-2} 5.10^{-2} $2.1.10^{-3}$ 10^{-3} 10^{-4} 10^{-4} 10^{-4}	$ \begin{array}{c} 10 \\ 1 \\ 6 \\ 3 \\ 4 \\ 5 \end{array} $	100 32 300 200 200 250	10^{-1} $2.5.10^{-3}$ $1.2.10^{-2}$ 2.10^{-4} 3.10^{-4} 6.10^{-4}	0 0 0 0 0 0	truss structure toy model toy model toy model hydraulics toy model	biased surrogate-based no surrogate, Lipschitz constant known no surrogate, g monotonic no surrogate, g monotonic no surrogate, g monotonic no surrogate, g monotonic	[54] [14] [68] [68] un. [68]

These results, and the previous considerations, lead us to the following conclusions, which in turn allow us to outline a research program.

First, we only scratch the surface on how to bound failure probabilities with powerful methodologies involving the choice of surrogates. Indeed, targeted probabilities in research papers remains too "high" with respect to typical orders of magnitude tied to severe industrial risks.

First, methodologies that are based on "biased" surrogates seem the most attractive to get non-asymptotic control of p, up to a given level of acceptability. While surrogates build on reduced bases usually come with bounds and could be ideal candidates, mixing constructive arguments proposed by [54] with other theoretical guarantees provided by concentration inequalities, as suggested by [39], could be a relevant way to produce such guarantees. Refining the construction of such surrogates through sequential DOE, and obtaining concentration inequalities on the basis of martingale-type arguments in the spirit of De La Peña exponential inequalities [12], seem a relevant avenue or research. Besides, ensuring the uniformity of first-order-stochastic dominance is probably not useful: we want to focus on some small subsets of Ω .

Second, it seems difficult, apart from some particular cases, to ensure strong geometrical assumptions directly on the behavior of $x \mapsto g(x)$. Perhaps the Lipschitz property seems more defensible than the monotonicity property in applications, but more probably, these two assumptions can be considered as true only for some part $X_{(1)}$ of the input Xconditionally to $X_{(2)} = X/X_{(1)}$. Provided the dimension of $X_{(1)}$ remains low, some small DOE could be used to provide conditional bounds on p (given $X_{(2)}$). Other arguments to get bounds from the model g itself could be inspired from the control techniques used on PDE solving by discretization (e.g., [40]. These techniques can provide hints to produce "cautious" surrogates, and offer an additional layer of guarantees for the end-user. Actually, the explicitly known or estimable geometric properties of these surrogates, such as their Lipschitz property and their monotonicity and convexity subdomains, make it possible to provide deterministic bounds on an estimator that is itself conservative and that would be produced from this surrogate, through the "biased surrogate" construction previously evoked.

In this sense, an interesting research idea could be to produce some additional criterion to the usual ones used for assessing the relevance of surrogates (e.g., Q^2) to ensure a form of conservative bias in interesting subsets of Ω . Such a work would fit with the recent development of a new criterion for Gaussian processes, interacting with the Q^2 , dedicated to ensure robust predictive properties of the surrogate [62].

In summary, the cases studied so far are usually very moderate in scale, and hardly reflect real cases (e.g. industrial, linked to critical systems), for which the need for a surrogate is particularly important. It therefore seems appropriate to launch a research program aimed at testing the scalability of the first methods based on first-order stochastic dominance, and to combine them with robust methods for calculating the bounds related to the geometry of these surrogates, through sequential explorations of the input space Ω .

Acknowledgments

The author thanks Bastien Bergère for having encoded and tested several ideas about monotonic Monte Carlo techniques, that have fed this review article.

A Appendix: Semi-adaptive sampling and computing bounds within staircase subspaces

As illustrated in Figure 2, boundary algorithms based on monotonicity require to explore nested staircase subspaces (defined after each sampling batch $A \to A'$ by the space between dominated subspaces $(\mathbb{U}^-(A'), \mathbb{U}^+(A')))$. Basically, any strategy for producing a statistical estimator, computing and improving the bounds must be based on uniform sampling in this "tortured" nondominated space. Rejection methods proposed in [18, 19] unfortunately require an increasing number of simulations as the as this space shrinks. A semi-adaptive Markov chain Monte Carlo method (MCMC) appears to be more efficient in dimension up to d = 15.

More formally, assume that a n-DOE of queries $x \mapsto g(x)$ has been chosen in $\Omega = [0, 1]^d$. This allows to define the two subspaces straightforwardly noted $\mathbb{U}_n^-, \mathbb{U}_n^+$, with $\mathbb{U}_n = \Omega/(\mathbb{U}_n^- \cup \mathbb{U}_n^+)$. Our goal is to sample a m-sample from the uniform distribution with density

$$f_n(x) = \frac{\mathbb{1}_{\{x \in \mathbb{U}_n\}}}{p^+ n - 1 - p^- n - 1}.$$

Doing so allows to estimate the new bounds $(p_n^-, p_n^+) = (\lambda(\mathbb{U}^-n), 1 - \lambda(\mathbb{U}_n^+))$ and select new candidate vectors for new queries $x \mapsto g(x)$. This can be done approximately using a MCMC technique, described beneath and illustrated on the successive graphs of Figures 6-7, such that the computation of bounds is consistent thanks to the ergodic theorem. To explore efficiently the staircase subspace using a traditional random walk instrumental distribution, a transformation $\psi : \Omega \to \mathbb{R}^d$ of each possible design vector $x \in \mathbb{U}_n$ is used, preserving the partial ordering $\leq_P: \psi = (\phi^{-1}, \dots, \phi^{-1})$ where ϕ the cumulative distribution function of the standard reduced Gaussian distribution.

More precisely, The semi-adaptiveness is defined in the following sense. Having constructed at a certain step n a Metropolis-Hastings mechanism, which has converged to the law $\mathcal{U}(\mathbb{U}_n)$, its trajectory to estimate the volumes of interest and simulate uniformly at steps $n + 1, \ldots, n + l$, for a fixed parameter l < n. A quasi-independent sample following $\mathcal{U}(\mathbb{U}_n)$ is obtained by batch sampling. Then, from this sample, one simulates uniformly in the sub-regions of interest by a simple accept-reject method.

1 Initialization

- **2** Given a nontrivial staircase subspace \mathbb{U}_0 , sample N values $\mathbf{X}^{(0)} \sim \mathcal{U}(\mathbb{U}_0)$.
- 3 While n < N do

 $\mathbf{5}$

6

4 1) Compute $\mathbf{Z} = \psi(X_{(n-l)})$ and the empirical covariance matrix of transformed past *l*-trajectories

$$\forall \ 1 \leq j,k \leq d \quad \hat{\Sigma}_{jk}^{(n-l)} = \frac{1}{N-1} \sum_{i=1}^{N} \left(Z_{ij} - \bar{Z}_j \right) \left(Z_{ik} - \bar{Z}_k \right)$$

then define the innovation $\varepsilon \sim \mathcal{N}_d(0, \alpha \hat{\Sigma}_{n-l}/d)$.

2) a) Build the Markov chain $\mathbf{X}^{(n)} = (X_1^{(n)}, ..., X_M^{(n)})$ from a Metropolis-Hastings mechanism on the transformed variable Z, based on the stationary distribution $\mathcal{U}(\mathbb{U}_n)$ and the transformed random walk proposal

$$\tilde{X} \stackrel{\mathcal{L}}{\sim} \psi^{-1} \left(\psi(X) + \varepsilon \right)$$

b) Using usual decorrelation tools, this leads to a nearly independent batch sample $\tilde{\mathbf{X}}^{(n)}$.

3) For each k = 1, ..., l, sample with rejection $X_{n+k} \sim \mathbb{U}_{n+k}$ from $\tilde{\mathbf{X}}^{(n)}$.

7 4) Estimate the volumes

$$\widehat{\lambda(\mathbb{U}_{n+k})} = \left(\frac{1}{M}\sum_{i=m+1}^{m+M}\mathbb{1}_{\{\mathbf{Y}_i^{(n)}\in\mathbb{U}_{n+k}\}}\right)\widehat{\lambda(\mathbb{U}_n)}$$
$$\widehat{\lambda(\mathbb{U}_{n+k}^+)} = \left(\frac{1}{M}\sum_{i=m+1}^{m+M}\mathbb{1}_{\{\mathbf{X}_i^{(n)}\in\hat{\mathbb{U}}_{n+k}^+\}}\right)\widehat{\lambda(\mathbb{U}_n)}$$

where, $\forall k \in [\![1, \ldots, l]\!], \forall i$:

$$\mathbb{1}_{\{\mathbf{Y}_{i}^{(n)} \in \mathbb{U}_{n+k}\}} = \left[\mathbb{1}_{\{\mathbf{Y}_{i}^{(n)} \geq X_{n+k}\}} \mathbb{1}_{\{\mathbf{X}_{n+k} \in \mathbb{U}-} + \mathbb{1}_{\mathbf{Y}_{i}^{(n)} \leq X_{n+k}\}} \mathbb{1}_{\{\mathbf{X}_{n+k} \in \mathbb{U}+\}}\right] \mathbb{1}_{\{\mathbf{Y}_{i}^{(n)} \in \mathbb{U}_{n+k-1}\}}$$

8 5) Update the bounds and update $n \leftarrow n+l$



Figure 6: Illustration in dimension 2 of the semi-adaptive MCMC method. From the previous chain which has converged (**A**), we obtain a quasi-iid sample following $\mathcal{U}(\mathbb{U}_4)$ (**B**), which we use to generate the next realisations $(\mathbf{X}_k)_{k=5,...,10}$ following the uniform distribution on the respective spaces $(\mathbb{U}_k)_{k=5,...,10}$ (**C**, **D**). The ergodic theorem is then applied to the entire trajectory of the chain to estimate the proportions of the sub-volumes $(\mathbb{U}_k)_{k=5,...,10}$ (**E**).



Figure 7: Illustration in dimension 2 of the semi-adaptive MCMC method. After step (**E**) on the previous figure, a point in the chain belonging to \mathbb{U}_{10} is chosen to initialise the next one (**F**). The Metropolis-Hastings random walk (**H**) is then evolved in the transformed space, taking Gaussian increments with a covariance equal to the empirical covariance of the previous chain (**G**). By inverse transformation we obtain a new chain converging to $\mathcal{U}(\mathbb{U}_{10})$ (**I**).

References

- Antoine Ajenjo, Emmanuel Ardillon, Vincent Chabridon, Bertrand Iooss, Scott Cogan, and Emeline Sadoulet-Reboul. An info-gap framework for robustness assessment of epistemic uncertainty models in hybrid structural reliability analysis. *Structural Safety*, 96:102196, 2022.
- [2] S. Amoukou. Trustworthy Machine Learning: Explainability and Distribution-Free Uncertainty Quantification. PhD thesis, Université Paris Saclay, 2023.
- [3] Siu-Kui Au and James L Beck. Estimation of small failure probabilities in high dimensions by subset simulation. *Probabilistic engineering mechanics*, 16(4):263–277, 2001.
- [4] Siu-Kui Au and James L Beck. Subset simulation and its application to seismic risk based on dynamic analysis. *Journal of engineering mechanics*, 129(8):901–917, 2003.
- [5] SK Au. Augmenting approximate solutions for consistent reliability analysis. Probabilistic Engineering Mechanics, 22(1):77–87, 2007.

- [6] Dario Azzimonti, David Ginsbourger, Clément Chevalier, Julien Bect, and Yann Richet. Adaptive design of experiments for conservative estimation of excursion sets. *Technometrics*, 63(1):13–26, 2021.
- [7] Julien Bect, François Bachoc, and David Ginsbourger. A supermartingale approach to gaussian process based sequential design of experiments. *Bernoulli*, 25:2883–2919, 2019.
- [8] Julien Bect, David Ginsbourger, Ling Li, Victor Picheny, and Emmanuel Vazquez. Sequential design of computer experiments for the estimation of a probability of failure. *Statistics and Computing*, 22:773–793, 2012.
- Julien Bect, Ling Li, and Emmanuel Vazquez. Bayesian subset simulation. SIAM/ASA Journal on Uncertainty Quantification, 5(1):762–786, 2017.
- [10] Yakov Ben-Haim. Info-gap decision theory: decisions under severe uncertainty. Elsevier, 2006.
- [11] Yakov Ben-Haim. Info-gap decision theory (ig). Decision making under deep uncertainty: From theory to practice, pages 93–115, 2019.
- [12] Bernard Bercu, Bernard Delyon, Emmanuel Rio, et al. Concentration inequalities for sums and martingales. Springer, 2015.
- [13] Bastien Bergère. Exploration de méthodes de monte-carlo accélérées sous contrainte de monotonie. *Master Thesis, Sorbonne Université'*, 2021.
- [14] Lucie Bernard, Albert Cohen, Arnaud Guyader, and Florent Malrieu. Recursive estimation of a failure probability for a lipschitz function. arXiv preprint arXiv:2107.13369, 2021.
- [15] Barron J Bichon, Michael S Eldred, Laura Painton Swiler, Sandaran Mahadevan, and John M McFarland. Efficient global reliability analysis for nonlinear implicit performance functions. AIAA journal, 46(10):2459–2468, 2008.
- [16] David Bolin and Finn Lindgren. Excursion and contour uncertainty regions for latent gaussian models. Journal of the Royal Statistical Society Series B: Statistical Methodology, 77(1):85–106, 2015.
- [17] J-M Bourinet, François Deheeger, and Maurice Lemaire. Assessing small failure probabilities by combined subset simulation and support vector machines. *Structural Safety*, 33(6):343–353, 2011.
- [18] Nicolas Bousquet. Accelerated monte carlo estimation of exceedance probabilities under monotonicity constraints. Annales de la Faculté des sciences de Toulouse: Mathématiques, 21(3):557–591, 2012.
- [19] Nicolas Bousquet, Thierry Klein, and Vincent Moutoussamy. Approximation of limit state surfaces in monotonic monte carlo settings, with applications to classification. SIAM/ASA Journal on Uncertainty Quantification, 6(1):1–33, 2018.
- [20] Charles-Edouard Bréhier, Tony Lelièvre, and Mathias Rousset. Analysis of adaptive multilevel splitting algorithms in an idealized case. *ESAIM: Probability and Statistics*, 19:361–394, 2015.
- [21] Karl Breitung. Sorm, design points, subset simulation, and markov chain monte carlo. ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering, 7(4):04021052, 2021.
- [22] Nader H Bshouty. Exact learning boolean functions via the monotone theory. Information and Computation, 123(1):146–153, 1995.

- [23] James Antonio Bucklew and J Bucklew. Introduction to rare event simulation, volume 5. Springer, 2004.
- [24] José-Ramón Cano, Pedro Antonio Gutiérrez, Bartosz Krawczyk, Michał Woźniak, and Salvador García. Monotonic classification: An overview on algorithms, performance measures and data sets. *Neurocomputing*, 341:168–182, 2019.
- [25] Frédéric Cérou, Pierre Del Moral, Teddy Furon, and Arnaud Guyader. Sequential monte carlo for rare event estimation. *Statistics and computing*, 22(3):795–808, 2012.
- [26] Frédéric Cérou, Arnaud Guyader, and Mathias Rousset. Adaptive multilevel splitting: Historical perspective and recent results. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 29(4), 2019.
- [27] Frédéric Cérou, Patrick Héas, and Mathias Rousset. Adaptive reduced multilevel splitting. arXiv preprint arXiv:2312.15256, 2023.
- [28] Kuan-Lin Chen, Harinath Garudadri, and Bhaskar D Rao. Improved bounds on neural complexity for representing piecewise linear functions. Advances in Neural Information Processing Systems, 35:7167–7180, 2022.
- [29] L.O. Chua and A.-C. Deng. Canonical piecewise-linear representation. IEEE Transactions on Circuits and Systems, 35(1):101–111, 1988.
- [30] Sergio Conti, Martin Rumpf, Ruüdiger Schultz, and Sascha Toölkes. Stochastic dominance constraints in elastic shape optimization. SIAM Journal on Control and Optimization, 56(4):3021–3034, 2018.
- [31] Sébastien Da Veiga and Amandine Marrel. Gaussian process modeling with inequality constraints. Annales de la Faculté des sciences de Toulouse: Mathématiques, 21(3):529– 555, 2012.
- [32] Sébastien Da Veiga and Amandine Marrel. Gaussian process regression with linear inequality constraints. *Reliability Engineering & System Safety*, 195:106732, 2020.
- [33] Andreas Damianou and Neil D Lawrence. Deep gaussian processes. In Artificial intelligence and statistics, pages 207–215. PMLR, 2013.
- [34] Marco de Angelis, Edoardo Patelli, and Michael Beer. Advanced line sampling for efficient robust reliability analysis. *Structural safety*, 52:170–182, 2015.
- [35] Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. A tutorial on the cross-entropy method. Annals of operations research, 134:19–67, 2005.
- [36] Etienne De Rocquigny. Structural reliability under monotony: Properties of form, simulation or response surface methods and a new class of monotonous reliability methods (mrm). *Structural Safety*, 31(5):363–374, 2009.
- [37] Ove Ditlevsen and Henrik O Madsen. Structural reliability methods, volume 178. Wiley New York, 1996.
- [38] Vincent Dubourg, Bruno Sudret, and Jean-Marc Bourinet. Reliability-based design optimization using kriging surrogates and subset simulation. *Structural and Multidisciplinary Optimization*, 44:673–690, 2011.
- [39] M. Ducoffe, S. Gerchinovitz, and J. Sen Gupta. A High-Probability Safety Guarantee for Shifted Neural Network Surrogates. *Proceedings of SafeAI 2020*, pages 74–82, 2020.
- [40] Daniel Elfverson, Fredrik Hellman, and Axel Målqvist. A multilevel monte carlo method for computing failure probabilities. SIAM/ASA Journal on Uncertainty Quantification, 4(1):312–330, 2016.

- [41] Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. Advances in Neural Information Processing Systems, 32, 2019.
- [42] Elias Fekhari, Bertrand Iooss, Joseph Muré, Luc Pronzato, and Maria-João Rendas. Model predictivity assessment: incremental test-set selection and accuracy evaluation. In *Convegno della Società Italiana di Statistica*, pages 315–347. Springer, 2021.
- [43] Joshua P French and Stephan R Sain. Spatio-temporal exceedance locations and confidence regions. Annals of Applied Statistics, 7:1421–1449, 2013.
- [44] Lea Friedli, David Ginsbourger, Arnaud Doucet, and Niklas Linde. An energy-based model approach to rare event probability estimation. arXiv:2310.04082, 2023.
- [45] N. Gayton, J.M. Bourinet, and M. Lemaire. Cq2rs: a new statistical approach to the response surface method for reliability analysis. *Structural Safety*, 25(1):99–121, 2003.
- [46] Peter W Glynn, Gerardo Rubino, and Bruno Tuffin. Robustness properties and confidence interval reliability issues. *Rare Event Simulation using Monte Carlo Methods*, pages 63–84, 2009.
- [47] O. Goldreich, S. Goldwassert, E. Lehman, and D. Ron. Testing monotonicity. In Proceedings 39th Annual Symposium on Foundations of Computer Science (Cat. No.98CB36280), pages 426–435, 1998.
- [48] Frank Grooteman. Adaptive radial-based importance sampling method for structural reliability. *Structural Safety*, 30(6):533–542, 2008.
- [49] Arnaud Guyader, Nicolas Hengartner, and Eric Matzner-Løber. Simulation and estimation of extreme quantiles and extreme probabilities. Applied Mathematics & Optimization, 64(2):171–196, 2011.
- [50] Abdul-Lateef Haji-Ali, Jonathan Spence, and Aretha Teckentrup. Adaptive multilevel monte carlo for probabilities. arXiv:2107.09148, 2021.
- [51] Boris Hanin and David Rolnick. Complexity of linear regions in deep networks. In International Conference on Machine Learning, pages 2596–2604. PMLR, 2019.
- [52] Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns. Advances in neural information processing systems, 32, 2019.
- [53] Jia Hao, Wenbin Ye, Liangyue Jia, Guoxin Wang, and Janet Allen. Building surrogate models for engineering problems by integrating limited simulation data and monotonic engineering knowledge. Advanced Engineering Informatics, 49:101342, 2021.
- [54] John D Jakeman, Drew P Kouri, and J Gabriel Huerta. Surrogate modeling for efficiently, accurately and conservatively estimating measures of risk. *Reliability Engineering & System Safety*, 221:108280, 2022.
- [55] P. L'Ecuyer, V. Demers, and B. Tuffin. Rare events, splitting, and quasi-Monte Carlo. ACM Transactions on Modeling and Computer Simulation (TOMACS), 17:Art. 9, 2007.
- [56] Nicolas Lelièvre, Pierre Beaurepaire, Cécile Mattrand, and Nicolas Gayton. Ak-mcsi: A kriging-based method to deal with small failure probabilities and time-consuming models. *Structural Safety*, 73:1–11, 2018.
- [57] Maurice Lemaire. Structural reliability. John Wiley & Sons, 2013.
- [58] Jing Li, Jinglai Li, and Dongbin Xiu. An efficient surrogate-based method for computing rare failure probability. *Journal of Computational Physics*, 230(24):8683– 8697, 2011.

- [59] Jing Li and Dongbin Xiu. Evaluation of failure probability via surrogate models. Journal of Computational Physics, 229(23):8966–8980, 2010.
- [60] Qui X Lieu, Khoa T Nguyen, Khanh D Dang, Seunghye Lee, Joowon Kang, and Jaehong Lee. An adaptive surrogate model to structural reliability analysis using deep neural network. *Expert Systems with Applications*, 189:116104, 2022.
- [61] Zhenzhou Lu, Shufang Song, Zhufeng Yue, and Jian Wang. Reliability sensitivity method by line sampling. *Structural Safety*, 30(6):517–532, 2008.
- [62] Amandine Marrel and Bertrand Iooss. Probabilistic surrogate modeling by gaussian process: A new estimation algorithm for more reliable prediction. HAL:cea-04322818, 2023.
- [63] Amandine Marrel, Bertrand Iooss, and Vincent Chabridon. The ICSCREAM methodology: Identification of penalizing configurations in computer experiments using screening and metamodel—applications in thermal hydraulics. *Nuclear Science and Engineering*, 196(3):301–321, 2022.
- [64] J. Morio, M. Balesdent, D. Jacquemart, and C. Vergé. A survey of rare event simulation methods for static input-output models. *Simulation Modelling Practice* and Theory, 49:287–304, 2014.
- [65] Jérôme Morio and Mathieu Balesdent. Estimation of rare event probabilities in complex aerospace and other systems: a practical approach. Woodhead publishing, 2015.
- [66] Maliki Moustapha, Stefano Marelli, and Bruno Sudret. Active learning for structural reliability: Survey, general framework and benchmark. *Structural Safety*, 96:102174, 2022.
- [67] Vincent Moutoussamy. Contributions à l'analyse de fiabilité structurale: prise en compte de contraintes de monotonie pour les modèles numériques. PhD thesis, Université Paul Sabatier-Toulouse III, 2015.
- [68] Vincent Moutoussamy, Nicolas Bousquet, Bertrand Iooss, Paul Rochet, Thierry Klein, and Fabrice Gamboa. Comparing conservative estimations of failure probabilities using sequential designs of experiments in monotone frameworks. In 11th International Conference on Structural Safety & Reliability (ICOSSAR), June 16-20, 2013, New York, 2013.
- [69] Scott Mueller and Judea Pearl. Monotonicity: Detection, refutation, and ramification. UCLA Technical report, 2023.
- [70] Sergei Ovchinnikov. Max-min representation of piecewise linear functions. arXiv preprint math/0009026, 2000.
- [71] Chirag Pabbaraju, Ezra Winston, and J Zico Kolter. Estimating lipschitz constants of monotone deep equilibrium models. In *International Conference on Learning Representations*, 2021.
- [72] Manolis Papadrakakis and Nikos D Lagaros. Reliability-based structural optimization using neural networks and monte carlo simulation. *Computer methods in applied* mechanics and engineering, 191(32):3491–3507, 2002.
- [73] Iason Papaioannou, Costas Papadimitriou, and Daniel Straub. Sequential importance sampling for structural reliability analysis. *Structural Safety*, 62:66–75, 2016.
- [74] Iason Papaioannou and Daniel Straub. Combination line sampling for structural reliability analysis. *Structural Safety*, 88:102025, 2021.

- [75] Victor Picheny, David Ginsbourger, Olivier Roustant, Raphael T Haftka, and Nam-Ho Kim. Adaptive designs of experiments for accurate approximation of a target region. *Journal of Mechanical Design*, 2010.
- [76] Angela Rita Provenzano, Daniele Trifirò, Alessio Datteo, Lorenzo Giada, Nicola Jean, Andrea Riciputi, G Le Pera, Maurizio Spadaccino, Luca Massaron, and Claudio Nordio. Machine learning approach for credit scoring. arXiv preprint arXiv:2008.01687, 2020.
- [77] Alfio Quarteroni, Andrea Manzoni, and Federico Negri. Reduced basis methods for partial differential equations: an introduction, volume 92. Springer, 2015.
- [78] Carl Edward Rasmussen and Hannes Nickisch. Gaussian processes for machine learning (gpml) toolbox. The Journal of Machine Learning Research, 11:3011–3015, 2010.
- [79] R Tyrrell Rockafellar and Stan Uryasev. The fundamental risk quadrangle in risk management, optimization and statistical estimation. Surveys in Operations Research and Management Science, 18(1-2):33-53, 2013.
- [80] Atin Roy and Subrata Chakraborty. Support vector machine in structural reliability analysis: A review. *Reliability Engineering & System Safety*, page 109126, 2023.
- [81] Thomas J Santner, Brian J Williams, William I Notz, and Brain J Williams. The design and analysis of computer experiments, volume 1. Springer, 2003.
- [82] Roland Schobi, Bruno Sudret, and Joe Wiart. Polynomial-chaos-based kriging. International Journal for Uncertainty Quantification, 5(2), 2015.
- [83] Timo Schorlepp, Shanyin Tong, Tobias Grafke, and Georg Stadler. Scalable methods for computing sharp extreme event probabilities in infinite-dimensional stochastic systems. *Statistics and Computing*, 2023.
- [84] Arnab Sharma and Heike Wehrheim. Higher income, larger loan? monotonicity testing of machine learning models. In Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis, pages 200–210, 2020.
- [85] Lorenzo Strigini and David Wright. Bounds on survival probability given mean probability of failure per demand; and the paradoxical advantages of uncertainty. *Reliability Engineering & System Safety*, 128:66–83, 2014.
- [86] R Sueur, A-L Popelin, and V Moutoussamy. Bounding and estimating failure probabilities in monotonic structural reliability. In *Safety, Reliability and Risk Analysis*, pages 1235–1243. CRC Press, 2013.
- [87] Shanyin Tong and Georg Stadler. Large deviation theory-based adaptive importance sampling for rare events in high dimensions. SIAM/ASA Journal on Uncertainty Quantification, 11(3):788–813, 2023.
- [88] Fulvio Tonon. Using random set theory to propagate epistemic uncertainty through a mechanical system. *Reliability Engineering & System Safety*, 85(1-3):169–181, 2004.
- [89] Elisabeth Ullmann. Rare event estimation with pde-based models: a tutorial. SIAM-UQ Conference, 2024.
- [90] Omar Valsson and Michele Parrinello. Variational approach to enhanced sampling and free energy calculations. *Physical review letters*, 113(9):090601, 2014.
- [91] Felipe AC Viana, Victor Picheny, and Raphael T Haftka. Using cross validation to design conservative surrogates. Aiaa Journal, 48(10):2286–2298, 2010.
- [92] Jun Wang and Zhiping Qiu. The reliability analysis of probabilistic and interval hybrid structural system. Applied Mathematical Modelling, 34(11):3648–3658, 2010.

- [93] Zebin Yang, Aijun Zhang, and Agus Sudjianto. Gami-net: An explainable neural network based on generalized additive models with structured interactions. *Pattern Recognition*, 120:108192, 2021.
- [94] Seungil You, David Ding, Kevin Canini, Jan Pfeifer, and Maya Gupta. Deep lattice networks and partial monotonic functions. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- [95] Xiukai Yuan, Shaolong Liu, Matthias Faes, Marcos A Valdebenito, and Michael Beer. An efficient importance sampling approach for reliability analysis of time-variant structures subject to time-dependent stochastic load. *Mechanical Systems and Signal Processing*, 159:107699, 2021.