



HAL
open science

Simulation conditionally to a subvariety and application to Bayesian optimization: a dichotomous approach

Frédéric Dambreville

► **To cite this version:**

Frédéric Dambreville. Simulation conditionally to a subvariety and application to Bayesian optimization: a dichotomous approach. SIMULTECH 2021, Jul 2021, Online, France. pp.178-203, 10.1007/978-3-031-23149-0_10 . hal-04004397

HAL Id: hal-04004397

<https://hal.science/hal-04004397v1>

Submitted on 24 Feb 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Simulation conditionally to a subvariety and application to Bayesian optimization: a dichotomous approach

Frédéric Dambreville¹[0000-0002-1460-4126]

ONERA/DTIS, Université Paris Saclay, F-91123 Palaiseau Cedex, France
frederic.dambreville@onera.fr

Abstract. This chapter concerns the issue of simulating a random vector conditionally to a subvariety. In practice, the subvariety is approximated by a thin shell surrounding it. These simulations are unusual cases of rare event simulation, which are made even more difficult with the dimension, the thinness of the shell and the nonlinearity of the subvariety. We propose a generic approach dedicated to the simulation of this type of rare event, which is dichotomous and inspired by interval analysis, with a particular effort to reduce the dimension curse. Non-incremental and incremental versions of the algorithm are designed in order to handle possible incremental constraint sets. Details on a multithread implementation are given especially in terms of data structure. These simulation methods are then used to propose a non-linear approach to Bayesian optimization. Examples of simulations and of Bayesian optimization are given in order to illustrate and compare the performances of the methods.

Keywords: Rare Event Simulation, Subvariety, Bayesian Optimization, Interval Analysis.

Notations

- \mathbb{R} is the set of reals and x, y, z are real variables,
- $[x] \triangleq [x^-, x^+]$, $[y]$, $[z]$ are notations for intervals,
- Bold notations $\mathbf{x} \triangleq \prod_{k=1}^n [x_k] = \prod_{k=1}^n [x_k^-, x_k^+]$ and $\mathbf{x}[\triangleq \prod_{k=1}^n [x_k^-, x_k^+[$ are used for boxes and half-open boxes respectively,
- $[\mathbb{R}] \triangleq \{[x^-, x^+] : [x^-, x^+] \subset \mathbb{R}\}$ is the set of interval subsets of \mathbb{R} .
 $[\mathbb{R}^n] \triangleq \{\prod_{k=1}^n [x_k] : \forall k, [x_k] \in [\mathbb{R}]\}$ is the set of box subsets of \mathbb{R}^n ,
- $\rho([x]) = x^+ - x^-$ and $\rho(\mathbf{x}) = \max_{1 \leq k \leq n} \rho([x_k])$ are the length of interval & box,
- $g : \mathbb{R}^k \rightarrow \mathbb{R}^j$ is a (multivariate) real function,
- $g(\mathbf{x}) \triangleq \{g(\mathbf{y}) : \mathbf{y} \in \mathbf{x}\}$. Set $g(\mathbf{x})$ is not necessarily a box,
- $[g] : [\mathbb{R}^k] \rightarrow [\mathbb{R}^j]$ is an interval function. Set $[g](\mathbf{x})$ is a box.

1 Introduction

This chapter is the long version of conference paper [3]. As a main topic, it addresses the issue of sampling a random vector defined on an n -dimension box

conditionally to a subvariety of this box. As a secondary topic and a main application of the simulation, the chapter proposes a non-linear algorithmic approach for Bayesian optimisation [10].

Thus, given a random vector \mathbf{X} of density $f_{\mathbf{X}}$ defined on \mathbb{R}^n , given box $[\mathbf{b}] \in [\mathbb{R}^n]$, *small* box $[\boldsymbol{\epsilon}] \in [\mathbb{R}^m]$ and map $g : [\mathbf{b}] \rightarrow \mathbb{R}^m$, our main objective is to sample conditional vector $[\mathbf{X} | g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$ ¹.

When space dimension n and constraint dimension m increase and length $\rho([\boldsymbol{\epsilon}])$ of $[\boldsymbol{\epsilon}]$ is small, then event $[g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$ has very low probability. We are dealing here with a particular case of rare event simulation. Moreover, as $[g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$ approximates a subvariety, it is foreseeable that conditional random vector $[\mathbf{X} | g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$ is essentially and “continuously” multimodal. Sketch figure 1 illustrates this well: to sample conditionally to the blue curve by a mixture of Gaussians, a large number of Gaussians are needed.

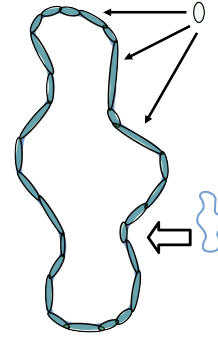


Fig. 1: Gaussian mixture

In survey [11], Morio, Balesdent et al. have evaluated the advantages and drawbacks of various rare event sampling methods. At this point, a comment may be done. In the literature, rare events are generally modelled by a function of risk being above an acceptable threshold. This formalism is quite general, but it suggests that the simulation of a rare event is tightly related to the maximization of a function. In practice, the maximizing set of a function is unimodal or somewhat multimodal. It is uncommon that this maximizing set is a subvariety. Not surprisingly, many methods evaluated in [11] are working well when the rare event is unimodal or moderately multimodal. In the case of conditional vector $[\mathbf{X} | g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$, these approaches do not work properly or need additional studies to take beneficially into account the subvariety structure.

This chapter promotes an alternative approach based on a dichotomous exploration in order to sample conditionally to a subvariety characterized by a known function. By design, this approach produces independent samples and avoids the phenomenon of sample impoverishment. These properties make it particularly well suited for an accurate approximation of a law conditionally to the subvariety. However, the method must fight the curse of dimensionality and we will be faced with two orthogonal dimensional problems: the exponential increase in sampling exploration and the degeneracy of particle weights. This paper presents an approach as well as generic and weakly parameterized algorithms, allowing good sampling performance for moderate dimensions (up to 11 for now) with balanced management of dimensional issues.

One main applicative purpose of this work is to contribute to the domain of Bayesian optimization from a nonlinear point of view. Bayesian optimization is among the key techniques for black-box optimization. In this kind of problems, one have to optimize fonctions which do not have usable formalization and for

¹ $[\mathbf{X} | g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$ is used as an abbreviation for $[\mathbf{X} | \mathbf{X} \in [\mathbf{b}] \ \& \ g(\mathbf{X}) \in [\boldsymbol{\epsilon}]]$.

which the request to an evaluator is generally particularly expensive. Accordingly, an inherent goal of such optimization is to save as much as possible on the number of candidate solutions to be evaluated. Approaches based on Gaussian surrogate models [8] have demonstrated the great interest of these Bayesian methods, but suffer from a drop in efficiency with the increase in dimension. Our contribution aims to address this dimensional difficulty by considering a known non-linear criterion depending on a model noise (a random vector) rather than a surrogate model of the criterion.

The chapter is divided in four main parts. Section 2 introduces concepts of interval analysis and derives some first ideas for a generic sampler. Section 3 deepens the intuitions introduced in section 2 and presents the working generic algorithms. Section 4 applies the sampling method for the design of a nonlinear algorithm for Bayesian optimization. Section 5 presents tests and analyses.

The algorithm already presented in conference paper [3] has slightly evolved and its implementation has been deeply revisited in terms of multithreading and data structure. In comparison to the conference paper, this chapter also describes a new incremental sampling algorithm, includes a new chapter dedicated to Bayesian optimization and presents rebuilt and supplemented tests and results.

2 Toward a generic sampling approach

This section presents some useful concepts and ideas for defining actually working sampling algorithms. Interval analysis is one of these tools. Interval analysis is a performing and accurate tool for dealing with constraints, and it provides a precise control on the approximation errors. Another interesting point is that intervals and boxes combine well with probability distributions.

2.1 Considerations about interval analysis

It is not our purpose to perform a good introduction on interval analysis [1, 7]. However, we refer to some key concepts, which are inspiring ideas for this work.

functions and operators. As a main ingredient throughout this chapter, it is reasonable to assume that for most common real functions g , it is possible to derive an interval function $[g]$ that satisfies properties:

- $g([\mathbf{x}]) \subset [g](\mathbf{x})$ for any $\mathbf{x} \in [\mathbb{R}^n]$,²
- $[g](\mathbf{x}) \subset [g](\mathbf{y})$ when $\mathbf{x} \subset \mathbf{y}$,
- If $\rho([\mathbf{x}])$ vanishes, then $\rho([g](\mathbf{x}))$ vanishes:

$$\rho(g([\mathbf{x}])) \xrightarrow{\rho([\mathbf{x}]) \rightarrow 0} 0 \tag{1}$$

²In particular, it is better if $[g]$ is minimal, that is $[g](\mathbf{x})$ is the minimal box containing $g(\mathbf{x})$ as subset for all $\mathbf{x} \in [\mathbb{R}^n]$.

These properties express that $[g]$ implies a bound on the error propagated by g , and this bound has good convergence behavior in regards to the error. Incidentally, we assume from now on that notation $[g]$ refers to an interval function, which verifies these good properties with respect to the function g .

Let us see how to build the interval functions on some common examples:

1. Reference functions, $g \in \{\ln, \exp, \sin, \cos, .^n, \dots\}$, are continuous onto \mathbb{R} , so that $g([x]) \in [\mathbb{R}]$ for all $[x] \in [\mathbb{R}]$. Then, it is optimal to set $[g]([x]) = g([x])$ for all $[x] \in [\mathbb{R}]$. As a consequence, the interval functions are easily optimally implementable for most reference functions. Here are some incomplete examples of definitions:

$$[\ln]([x]) \triangleq [\ln(x^-), \ln(x^+)] \quad \text{and} \quad [x]^n \triangleq \begin{cases} [x_-^n, x_+^n] & \text{for } n > 0, \\ [x_+^n, x_-^n] & \text{for } n < 0 \text{ \& } 0 \notin [x], \end{cases} \quad (2)$$

$$[\cos][x] \triangleq \begin{cases} [\cos(x^-), \cos(x^+)] & \text{for } [x] \subset [-\pi, 0] \\ [\cos(x^+), \cos(x^-)] & \text{for } [x] \subset [0, \pi] \\ [\min(\cos(x^-), \cos(x^+)), 1] & \text{for } 0 \in [x] \subset [-\pi, \pi] \\ \text{etc.} \end{cases} \quad (3)$$

2. Minimal interval functions for classical operators $+$, \cdot , $-$, $/$ are also easily defined. For example:

$$[x] + [y] \triangleq [x^- + y^-, x^+ + y^+] \quad \text{and} \quad [x] - [y] \triangleq [x^- - y^+, x^+ - y^-]. \quad (4)$$

3. Function g defined by $g(\theta) = (\cos(\theta), \sin(\theta))$ is an example, which is such that $g([\theta]) \notin [\mathbb{R}^2]$. One would rather define $[g]([\theta]) = [\cos]([\theta]) \times [\sin]([\theta])$ (\times is the Cartesian product) which is a strict superset of $g([\theta])$ in general. By the way, this construction is an illustration on how the interval functions of reference are used to define complex interval functions straightforwardly.

There is no uniqueness in the construction of $[g]$, unless it is chosen minimal. Unfortunately, constructing a minimal $[g]$ is not always easy or automatic. Let us consider the case of function $g : \theta \mapsto \cos^2 \theta + \sin^2 \theta$. Then, there are two obvious definitions for $[g]$:

1. By using the reference functions \cos , \sin , $.^2$ and $+$ one may derive:

$$[g]([\theta]) = ([\cos]([\theta]))^2 + ([\sin]([\theta]))^2. \quad (5)$$

For example, we compute $[g]([0, \frac{\pi}{2}]) = [0, 2]$ which is a bad error bound on theoretical value 1. Now, we also compute $[g]([-\frac{1}{10}, \frac{1}{10}]) \simeq [0.99, 1.01]$ which is a tight error bound on 1. This example holds confirmation that $[g]([\theta])$ has a good behavior for small boxes $[\theta]$.

2. By noticing that $\cos^2 \theta + \sin^2 \theta = 1$, it is optimal to define $[g]([\theta]) = [1, 1]$.

Although approach 2 gives the best solution, in practice approach 1 is preferred since it is generic, it is based on already implemented functions of reference, and it provides a way to construct automatically $[g]$ without any specific knowledge.

Subpaving and set inversion. As $[g]$ implies a bound on errors propagated by g with good convergence behavior, it may be used combined with a dichotomous process to produce a subpaving which efficiently approximates a set inversion $g^{-1}([\mathbf{y}])$. The example on figure 2 is kindly given by professor Jaulin and is also taken from [4]. It shows a resulting subpaving which approximates a set inversion. The decomposition is clearly dichotomous.

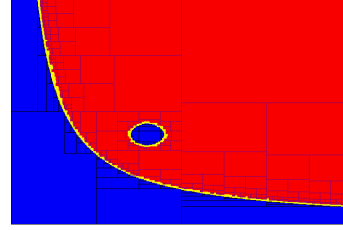


Fig. 2: Subpaving

A bisection process is iterated starting from main box $[\mathbf{b}]$; at each iteration, sub-boxes $[\mathbf{x}]$ are tested against the constraint $g([\mathbf{x}]) \subset [\mathbf{y}]$. Three cases arise:

- case a:** $g([\mathbf{x}]) \subset [\mathbf{y}]$, then $[\mathbf{x}]$ is among *red* boxes, which constitute a *subpaving* of $g^{-1}([\mathbf{y}])$.
- case b:** $g([\mathbf{x}]) \cap [\mathbf{y}] = \emptyset$, then $[\mathbf{x}]$ is among *blue* boxes, which constitute a *subpaving* of $\mathbb{R}^n \setminus g^{-1}([\mathbf{y}])$.
- case c:** Otherwise, bisection has to be repeated on $[\mathbf{x}]$ until sufficient convergence (*yellow* color).

Property (1) plays a key role in the decomposition process, ensuring that case **a** or case **b** are finally achieved when sub-boxes $[\mathbf{x}]$ are sufficiently small and sufficiently far from the frontier of set $g^{-1}([\mathbf{y}])$.

2.2 Naive dichotomous approach for sampling

Let μ be Borel measure on \mathbb{R}^n . It is given from now on:

- a random vector \mathbf{X} on \mathbb{R}^n characterized by a *bounded* density $f_{\mathbf{X}}$. Cumulative distribution function $F_{\mathbf{X}}$ of \mathbf{X} , defined by $F_{\mathbf{X}}(\mathbf{x}) = P(\mathbf{X} \leq \mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$, is assumed to be easily computable,
- a box $[\mathbf{b}] \in [\mathbb{R}^n]$ and a small box $[\boldsymbol{\epsilon}] \in [\mathbb{R}^m]$,
- a continuous map $g : [\mathbf{b}] \rightarrow \mathbb{R}^m$ built of functions and operators of reference,
- $[g]$ derived from g and related interval functions and operators of reference.

Sampling within boxes. We point out that it is easy to compute $P(\mathbf{X} \in [\mathbf{y}])$ or to sample $[\mathbf{X} | \mathbf{X} \in [\mathbf{y}]]$ when $F_{\mathbf{X}}$ is available, especially when the components of \mathbf{X} are jointly independent. These results are well known, and details are given in [4]. Thus, these features are taken for granted in this paper.

Sampling by means of a subpaving. Assume that set $g^{-1}([\boldsymbol{\epsilon}])$ has been approximated (by excess) by a subpaving. Thus, there is $\mathbb{P} \subset [\mathbb{R}^n]$ such that:

- For all $[\mathbf{x}], [\mathbf{y}] \in \mathbb{P}$, boxes $[\mathbf{x}[$ and $[\mathbf{y}[$ are disjoint,
- $g^{-1}([\boldsymbol{\epsilon}]) \subseteq \bigcup_{[\mathbf{x}] \in \mathbb{P}} [\mathbf{x}$ and $\sqcup_{\mathbb{P}} \simeq g^{-1}([\boldsymbol{\epsilon}])$, where $\sqcup_{\mathbb{P}} \triangleq \bigsqcup_{[\mathbf{x}] \in \mathbb{P}} [\mathbf{x}[$.

Set \mathbb{P} is typically composed from boxes of case **a** and case **c** (red and yellow colors) after a dichotomous subpaving construction of the set inversion.

The quality of the approximation may be quantified by set measures:

$$\alpha_{\mathbb{P}} = \mu\left(\sqcup_{\mathbb{P}} \setminus g^{-1}([\epsilon])\right) = \sum_{[\mathbf{x}] \in \mathbb{P}} \mu([\mathbf{x}]) - \mu(g^{-1}([\epsilon])). \quad (6)$$

Smaller is $\alpha_{\mathbb{P}}$, better is the approximation. Interval based set inversions are able to reach arbitrary precision for small dimensions (2 or 3 typically).

Now, for all $\mathbf{y} \in \mathbb{R}^n$, it happens that:

$$\begin{aligned} f_{\mathbf{X}|\mathbf{X} \in \sqcup_{\mathbb{P}}}(\mathbf{y}) &= \frac{f_{\mathbf{X}}(\mathbf{y})\delta_{\mathbf{y} \in \sqcup_{\mathbb{P}}}}{P(\mathbf{X} \in \sqcup_{\mathbb{P}})} & (7) \\ &= \frac{f_{\mathbf{X}}(\mathbf{y}) \sum_{[\mathbf{x}] \in \mathbb{P}} \delta_{\mathbf{y} \in [\mathbf{x}]} P(\mathbf{X} \in [\mathbf{x}])}{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}])} = \frac{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}]) f_{\mathbf{X}|\mathbf{X} \in [\mathbf{x}]}(\mathbf{y})}{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}])} \\ &= \sum_{[\mathbf{x}] \in \mathbb{P}} \frac{P(\mathbf{X} \in [\mathbf{x}])}{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}])} f_{\mathbf{X}|\mathbf{X} \in [\mathbf{x}]}(\mathbf{y}), & (8) \end{aligned}$$

where $\delta_{\text{true}} = 1$ and $\delta_{\text{false}} = 0$ else. Since $f_{\mathbf{X}}$ is bounded, equation (7) implies:

$$[\mathbf{X} | \mathbf{X} \in \sqcup_{\mathbb{P}}] \xrightarrow[\alpha_{\mathbb{P}} \rightarrow 0]{\mathcal{L}} [\mathbf{X} | \mathbf{X} \in g^{-1}([\epsilon])] \quad (9)$$

Now, equation (8) shows clearly that $f_{\mathbf{X}|\mathbf{X} \in \sqcup_{\mathbb{P}}}$ may be sampled by applying two steps: first sample a box $[\mathbf{x}] \in \mathbb{P}$ according to the discrete probability $\frac{P(\mathbf{X} \in [\mathbf{x}])}{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}])}$, then sample \mathbf{y} by the conditional law $f_{\mathbf{X}|\mathbf{X} \in [\mathbf{x}]}$. At last, we have

here an efficient method for sampling $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$ (algorithm 1).

The approach is efficient on conditional events like $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$. But this application of the interval-based inversion is only applicable to rather small dimensions. Taking inspiration of this preliminary approach, we address now the sampling problem in higher dimensions.

```

Function Sampling  $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$ 
  input :  $\alpha, g, N$    output:  $\mathbf{y}_{1:N}$ 
  Build subpaving  $\mathbb{P}$  such that  $\alpha_{\mathbb{P}} < \alpha$ 
  for  $k \leftarrow 1$  to  $N$  do
    Select  $[\mathbf{x}] \in \mathbb{P}$  with proba.  $\frac{P(\mathbf{X} \in [\mathbf{x}])}{\sum_{[\mathbf{x}] \in \mathbb{P}} P(\mathbf{X} \in [\mathbf{x}])}$ 
    Build  $\mathbf{y}_k$  by sampling  $[\mathbf{X} | \mathbf{X} \in [\mathbf{x}]]$ 
  end

```

Algorithm 1: Based on a subpaving.

Toward an improved approach. A key point of algorithm 1 is to be able to sample a box $[\mathbf{x}]$ of a subpaving of $g^{-1}([\epsilon])$. It is noticeable that an entire build of the subpaving is not needed here. Indeed, if we were able to construct a box $[\mathbf{x}]$ of the subpaving *on demand*, together with its *relative weight* within the subpaving, then we would be able to build sample \mathbf{y} .

Therefore, it may be opportune to merge the sampling process with the dichotomous construction of the subpaving itself. Now, a main ingredient of a dichotomous approach is also how the algorithm divides and conquers. In general, bisections are often used in dichotomous processes, as there is a guaranty of exponential volume decrease of the search area. Our approach is less constrained, so as to better tune the exploration strategy.

```

1 Function Sampling [ $\mathbf{X} | g(\mathbf{X}) \in [\epsilon]$ ]
   input :  $r, g, \omega, N$    output:  $(\mathbf{y}_k, w_k)_{1:N}$ 
2   for  $k \leftarrow 1$  to  $N$  do
3      $([\mathbf{x}_0], \pi_k, j) \leftarrow ([\mathbf{b}], 1, 0)$ 
4     while  $\rho([\mathbf{x}_j]) > r$  and  $[g]([\mathbf{x}_j]) \not\subset [\epsilon]$  do
5        $([\mathbf{l}_{j+1}], [\mathbf{r}_{j+1}]) \leftarrow \text{Cut}([\mathbf{x}_j])$ 
6        $j \leftarrow j + 1$ 
7        $[\mathbf{x}_j] \leftarrow \text{Bern}([\mathbf{l}_j], \omega_{[\mathbf{l}_j]}), ([\mathbf{r}_j], \omega_{[\mathbf{r}_j]})$ 
8        $\pi_k \leftarrow \frac{\omega_{[\mathbf{x}_j]}}{\omega_{[\mathbf{l}_j]} + \omega_{[\mathbf{r}_j]}} \pi_k$ 
9     end
10     $w_k \leftarrow P(\mathbf{X} \in [\mathbf{x}_j]) / \pi_k$ 
11    Build  $\mathbf{y}_k$  by sampling  $[\mathbf{X} | \mathbf{X} \in [\mathbf{x}_j]]$ 
12  end
13 end

```

Algorithm 2: Based on a weighting function.

Here, we speak in terms of cuts, which are more general, *being implied that an appropriate management of the box length is made in order to ensure the convergence*. In this paper, a cut is defined as follows:

- A cut of box $[\mathbf{x}] \in [\mathbb{R}^n]$ is a pair $([\mathbf{l}], [\mathbf{r}]) \in [\mathbb{R}^n]^2$ such that:

$$[\mathbf{l}] \cap [\mathbf{r}] = \emptyset \quad \text{and} \quad [\mathbf{l}] \sqcup [\mathbf{r}] = [\mathbf{x}], \quad (10)$$

- A bisection is a cut $([\mathbf{l}], [\mathbf{r}])$ such that $[\mathbf{l}]$ and $[\mathbf{r}]$ are same-sized.

In order to drive the dichotomous sampling process, we assume that a predictive weighting function is available:

$$\begin{cases} \omega_{[\mathbf{x}]} = 0 & \text{if } [g]([\mathbf{x}]) \cap [\epsilon] = \emptyset \\ \omega_{[\mathbf{x}]} \simeq P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon]) & \text{otherwise.} \end{cases} \quad (11)$$

Algorithm 2 implicitly builds a partial subpaving, and produces at the same time a weighted particle cloud as a result of the sampling of $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$. The algorithm iterates (**for** loop) the same sampling process, that is the following successive steps, until $[\mathbf{x}_j]$ is sufficiently small (*i.e.* $\rho([\mathbf{x}_j]) \leq r$) or is inside an implied suppaving (*i.e.* $[g]([\mathbf{x}_j]) \subset [\epsilon]$)³:

- Build a cut of $[\mathbf{x}_j]$ by means of function $\text{Cut}([\mathbf{x}_j])$. This function is designed so as to ensure that $\rho([\mathbf{x}_j])$ vanishes,
- Select randomly one box of cut $([\mathbf{l}_j], [\mathbf{r}_j])$ in proportion to their weight, by mean of Bernoulli process $\text{Bern}([\mathbf{l}_j], \omega_{[\mathbf{l}_j]}), ([\mathbf{r}_j], \omega_{[\mathbf{r}_j]})$,
- Update π_k which computes the processed probability of $[\mathbf{x}_j]$ in regards to the Bernoulli sequence.

³Recall that $[g]([\mathbf{x}_j])$ is easily computable while $g([\mathbf{x}_j])$ is not.

Assume that J is the last value reached by parameter j after the **while** loop. Then, the corrected weight $w_k = \frac{1}{\pi_k} P(\mathbf{X} \in [\mathbf{x}_J])$ is computed for $[\mathbf{x}_J]$ and for \mathbf{y}_k , and \mathbf{y}_k is sampled from $[\mathbf{x}_J]$.

Notice that $\omega_{[\mathbf{x}]} = 0$ when $[g]([\mathbf{x}]) \cap [\epsilon] = \emptyset$, so that boxes $[\mathbf{x}_J]$ are necessary within a subpaving of $g^{-1}([\epsilon])$ or its border, thanks to the Bernoulli process. Combined with loop constraint $\rho([\mathbf{x}_j]) > r$ **and** $[g]([\mathbf{x}_j]) \not\subset [\epsilon]$, it follows that a subpaving of $g^{-1}([\epsilon])$ (or its border) is implicitly and partially built during the sampling process.

When $[g]([\mathbf{x}_J]) \subset [\epsilon]$, we have $w_k = \frac{P(\mathbf{X} \in [\mathbf{x}_J])}{\pi_k}$ where π_k evaluates the processed probability for $[\mathbf{x}_J]$. As a result, the weighted particles (\mathbf{y}_k, w_k) provide an unbiased estimation of $f_{\mathbf{X}|g(\mathbf{X}) \in [\epsilon]}$ in a subpaving of $g^{-1}([\epsilon])$. It is not the same at the border of $g^{-1}([\epsilon])$, but this case is neglected. However, the sampler is not at all efficient when considering its variance.

Assume $\omega_{[\mathbf{x}]} = P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$, a case which works perfectly. In this ideal case, the weight along a **while** loop is computed by:

$$\frac{\omega_{[\mathbf{x}_j]}}{\omega_{[\mathbf{l}_j]} + \omega_{[\mathbf{r}_j]}} = \frac{P(\mathbf{X} \in [\mathbf{x}_j] \& g(\mathbf{X}) \in [\epsilon])}{P(\mathbf{X} \in [\mathbf{l}_j] \cup [\mathbf{r}_j] \& g(\mathbf{X}) \in [\epsilon])} = \frac{P(\mathbf{X} \in [\mathbf{x}_j] \& g(\mathbf{X}) \in [\epsilon])}{P(\mathbf{X} \in [\mathbf{x}_{j-1}] \& g(\mathbf{X}) \in [\epsilon])}, \quad (12)$$

and then:

$$\pi_k = \prod_{j=1}^J \frac{\omega_{[\mathbf{x}_j]}}{\omega_{[\mathbf{l}_j]} + \omega_{[\mathbf{r}_j]}} = \frac{P(\mathbf{X} \in [\mathbf{x}_J] \& g(\mathbf{X}) \in [\epsilon])}{P(\mathbf{X} \in [\mathbf{b}] \& g(\mathbf{X}) \in [\epsilon])}. \quad (13)$$

Three cases potentially arise:

- $[g]([\mathbf{x}_J]) \subset [\epsilon]$, *i.e.* $[\mathbf{x}_J]$ is in implied subpaving. Since $g([\mathbf{x}_J]) \subset [g]([\mathbf{x}_J])$, it comes $P(\mathbf{X} \in [\mathbf{x}_J] \& g(\mathbf{X}) \in [\epsilon]) = P(\mathbf{X} \in [\mathbf{x}_J])$. Then:

$$\begin{aligned} w_k &= \frac{P(\mathbf{X} \in [\mathbf{x}_J])}{\frac{P(\mathbf{X} \in [\mathbf{x}_J] \& g(\mathbf{X}) \in [\epsilon])}{P(\mathbf{X} \in [\mathbf{b}] \& g(\mathbf{X}) \in [\epsilon])}} = \frac{P(\mathbf{X} \in [\mathbf{x}_J])}{\frac{P(\mathbf{X} \in [\mathbf{x}_J])}{P(\mathbf{X} \in [\mathbf{b}] \& g(\mathbf{X}) \in [\epsilon])}} = P(\mathbf{X} \in [\mathbf{b}] \& g(\mathbf{X}) \in [\epsilon]) \\ &= P(g(\mathbf{X}) \in [\epsilon]) \end{aligned} \quad (14)$$

- $[g]([\mathbf{x}_J]) \cap [\epsilon] \neq \emptyset$ but $[g]([\mathbf{x}_J]) \not\subset [\epsilon]$, *i.e.* $[\mathbf{x}_J]$ is within the border of the implied subpaving. These cases are negligible for small precision r .
- $[g]([\mathbf{x}_J]) \cap [\epsilon] = \emptyset$, *i.e.* $[\mathbf{x}_J]$ is outside the implied subpaving and its border. This case is simply impossible from the Bernoulli process.

Equation (14) shows that the sampling process results in a cloud of same-weight particles over the implied subpaving. Border cases are negligible. Here we have a sampler of $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$ with the best variance performance in regards to the number of particles. But hypothesis $\omega_{[\mathbf{x}]} = P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$ is necessary. Of course such exact weighting function is almost never available.

Why Does It Generally Not Work? When $\omega_{[\mathbf{x}]} \neq P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$, the accumulated error will explode with the dimension, which will result in dramatically uneven weights on the particles. The resulting weighted particles cloud is then useless for practical applications.

3 Generic dichotomous approaches for sampling

Algorithms 1 and 2 illustrate the two main dimensional issues, that we have to deal with. These approaches are complementary:

- By building a complete subpaving of $g^{-1}([\epsilon])$, algorithm 1 makes possible a direct sampling of $[\mathbf{X} | g(\mathbf{X}) \in [\epsilon]]$, and incidently an accurate computation of $P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$. However, this construction of a complete subpaving is only possible for small dimensions.
- Algorithm 2 avoids the construction of a complete subpaving. Instead, it builds the boxes of an implied subpaving on demand throughout the sampling iteration. However, the algorithm is inefficient unless the predictive weighting function $\omega_{[\mathbf{x}]}$ is a good approximation of $P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$. This condition is not accessible in general.

We propose now an intermediate approach which:

- keeps history of the subpaving construction throughout the sampling process,
- use this history to build an improved estimate of $P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$.

By these tricks, it is expected that the sampling precision will increase with the number of samples. In order to avoid useless exploration, we also truncate the dichotomous process on the basis of some predictive assessment of final weight w_k . Thus, the algorithm tends to favor breadth search instead of depth search at the early stages of the sampling process.

3.1 Some containment of the curse of dimension

From now on, it is assumed that:

$$0 \leq \omega_{[\mathbf{x}]} \leq P(X \in [\mathbf{x}]) , \quad (15)$$

and that:

$$\omega_{[\mathbf{x}]} = \begin{cases} 0 & \text{if } [g]([\mathbf{x}]) \cap [\epsilon] = \emptyset , \\ P(\mathbf{X} \in [\mathbf{x}]) & \text{if } [g]([\mathbf{x}]) \subset [\epsilon] . \end{cases} \quad (16)$$

Algorithm 3 is an evolution of algorithm 2. In addition, it builds an history of cuts, stored in map `cuts`, and computes dynamically from this history an improved weighting function, stored in map `omg`.

The lines of this algorithm are colored in black, blue or dark blue. Black lines are inherited from algorithm 2. Blue lines are new additions to the previous algorithm. Dark blue lines (4, 19, 23 and 2 partially) correspond to the **for** loop of algorithm 2 rewritten as a **while** loop:

```
k ← 0 while k < N do ... k ← k + 1 ... end
```

Variable `cuts` is a dictionary and is used to register the history of computed cuts. At start, `cuts` is defined empty (line 2). For a given box $[\mathbf{x}_j]$, the cut on $[\mathbf{x}_j]$ is computed only once, if it is computed, by line 8:

```
ifundef cuts([\mathbf{x}_j]) ← Cut([\mathbf{x}_j])
```

```

1 Function Sampling [ $\mathbf{X} | g(\mathbf{X}) \in [\epsilon]$ ]
   input :  $\sigma, r, g, \omega, N$  output:  $(\mathbf{y}_k, w_k)_{1:N}$ 
2    $(\text{cuts}, \text{omg}, k) \leftarrow (\emptyset, \emptyset, 0)$ 
3    $\text{omg}([\mathbf{b}]) \leftarrow \omega_{[\mathbf{b}]}$ 
4   while  $k < N$  do
5      $([\mathbf{x}_0], \pi_k, j) \leftarrow ([\mathbf{b}], 1, 0)$ 
6     while  $\rho([\mathbf{x}_j]) > r$  and  $[g]([\mathbf{x}_j]) \not\subset [\epsilon]$  do
7       if  $\left| \log_2 \left( \frac{\text{omg}([\mathbf{b}])}{\text{omg}([\mathbf{x}_j])} \pi_k \right) \right| > \sigma$  goto 20
8       ifundef  $\text{cuts}([\mathbf{x}_j]) \leftarrow \text{Cut}([\mathbf{x}_j])$ 
9        $([l_{j+1}], [r_{j+1}]) \leftarrow \text{cuts}([\mathbf{x}_j])$ 
10       $j \leftarrow j + 1$ 
11      ifundef  $\text{omg}([r_j]) \leftarrow \omega_{[r_j]}$ 
12      ifundef  $\text{omg}([l_j]) \leftarrow \omega_{[l_j]}$ 
13       $(\nu_{[l_j]}, \nu_{[r_j]}) \leftarrow (\text{omg}([l_j]), \text{omg}([r_j]))$ 
14       $[\mathbf{x}_j] \leftarrow \text{Bern}([l_j], \nu_{[l_j]}, ([r_j], \nu_{[r_j]}))$ 
15       $\pi_k \leftarrow \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[l_j]} + \nu_{[r_j]}} \pi_k$ 
16    end
17     $w_k \leftarrow P(\mathbf{X} \in [\mathbf{x}_j]) / \pi_k$ 
18    Build  $\mathbf{y}_k$  by sampling  $[\mathbf{X} | \mathbf{X} \in [\mathbf{x}_j]]$ 
19     $k \leftarrow k + 1$ 
20    for  $i \leftarrow j$  to 1 do
21       $\text{omg}([\mathbf{x}_{i-1}]) \leftarrow \text{omg}([l_i]) + \text{omg}([r_i])$ 
22    end
23  end
24 end

```

Algorithm 3: Based on cuts history.

Keyword **ifundef** tests if $\text{cuts}([\mathbf{x}_j])$ is defined; if still undefined, then $\text{cuts}([\mathbf{x}_j])$ is set to $\text{Cut}([\mathbf{x}_j])$.

Variable **omg** is a dictionary which records the predictive weighting function and its possible updates, when needed. At start, **omg** is only defined for $[\mathbf{b}]$ and is set to $\omega_{[\mathbf{b}]}$ (lines 2 and 3). Variable $\text{omg}([r_j])$ is set to $\omega_{[r_j]}$, if it has not been initialized yet (line 11). The same is done for variable $\text{omg}([l_j])$ at line 12. When the cuts sequence is done (second **while**), then the weighting function is updated by the **for** loop (lines 20, 21, 22). This ensures the computation of $\text{omg}([\mathbf{x}])$ as the sum of the weights $\text{omg}([\mathbf{z}])$ of the leaves $[\mathbf{z}]$ of the cuts tree rooted on $[\mathbf{x}]$. Then, property (16) ensures that $\text{omg}([\mathbf{x}])$ gets closer to $P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\epsilon])$ when the cuts tree rooted on $[\mathbf{x}]$ gets more refined.

Algorithm 3 is similar to algorithm 2, except that:

- cut $([l_{j+1}], [r_{j+1}])$ is recovered from the history, when it is possible (line 9),
- box selection is done by means of $(\nu_{[l_j]}, \nu_{[r_j]}) \triangleq (\text{omg}([l_j]), \text{omg}([r_j]))$.

There is an interesting property here. Assume that J is the last value reached by j and that $J' < J$ is such that $\text{omg}([l_j])$ and $\text{omg}([r_j])$ are already defined for

all $1 \leq j \leq J'$. Weighting functions are updated in these cases. Then, it comes for all $1 \leq j \leq J'$ that:

$$\text{omg}([\mathbf{x}_{j-1}]) = \text{omg}([\mathbf{l}_j]) + \text{omg}([\mathbf{r}_j]).$$

The computation of π_k is then simplified:

$$\begin{aligned} \pi_k &= \prod_{j=1}^{J'} \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[\mathbf{l}_j]} + \nu_{[\mathbf{r}_j]}} \prod_{j=J'+1}^J \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[\mathbf{l}_j]} + \nu_{[\mathbf{r}_j]}} = \prod_{j=1}^{J'} \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[\mathbf{x}_{j-1}]}} \prod_{j=J'+1}^J \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[\mathbf{l}_j]} + \nu_{[\mathbf{r}_j]}} \\ &= \frac{\nu_{[\mathbf{x}_{J'}]}}{\nu_{[\mathbf{b}]}} \prod_{j=J'+1}^J \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[\mathbf{l}_j]} + \nu_{[\mathbf{r}_j]}}. \end{aligned} \quad (17)$$

Thus, the error on π_k grows exponentially only within the newly explored cuts, that is here from $J' + 1$ to J . This is a reason for setting a certain restriction on the depth-oriented aspect of this sampling process. Another good reason is to prevent degenerate particle weights, w_k . Algorithm 3 thus implements some code (line 7) for testing the degeneracy of π_k and eventually restarting the sampling loop (second **while**):

$$\text{if } \left| \log_2 \left(\frac{\text{omg}([\mathbf{b}])}{\text{omg}([\mathbf{x}_j])} \pi_k \right) \right| > \sigma \text{ goto } 20$$

This code tests the logarithmic distance between the weight of $[\mathbf{x}_j]$, $\text{omg}([\mathbf{x}_j])$, and the weight resulting from the sampling process, $\text{omg}([\mathbf{b}]) \pi_k$. If it is higher than σ , then the loop is stopped by going to line 20. By doing that, the incrementation of k is skipped, so that the sampling loop is restarted for the same indice k . However, the update of variable omg is done, and of course, the history of cuts stays incremented. So, although the sampling loop has been interrupted in this case, the sampling structure has been upgraded. This results in an adaptive process which will balance depth and breadth explorations when running the sampling. Breadth exploration is favored on the first sampling iterations, but the tendency becomes inverted after several samples.

3.2 Incremental algorithm

In section 4, we present an application of the conditional simulation for Bayesian optimization. In this applicative context, we have to successively simulate a random vector conditionally to incremental constraints.

More precisely, let $U \geq 1$ and let $[\epsilon_u] \in [\mathbb{R}^{m_u}]$ and $g_u : [\mathbf{b}] \rightarrow \mathbb{R}^{m_u}$ be defined for $1 \leq u \leq U$. For $1 \leq v \leq U$ are defined:

$$[\epsilon_{1:v}] = \prod_{u=1}^v [\epsilon_u], \quad (18)$$

and $g_{1:v} : [\mathbf{b}] \rightarrow \prod_{u=1}^v \mathbb{R}^{m_u}$ by:

$$g_{1:v}(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_v(\mathbf{x})) \quad \text{for all } \mathbf{x} \in [\mathbf{b}]. \quad (19)$$

```

1 Function Sampling [ $\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]$ ]
2   input :  $\sigma, r, g, \omega, N$    output:  $(\mathbf{y}_k, w_k)_{1:N}$ 
3    $(\text{cuts}, \text{omg}, \text{lev}) \leftarrow (\emptyset, \emptyset, \emptyset)$ 
4   for  $v \leftarrow 1$  to  $U$  do
5      $k \leftarrow 0$ 
6     ifundef  $(\text{omg}([\mathbf{b}]), \text{lev}([\mathbf{b}])) \leftarrow (\omega_{[\mathbf{b}]}^{1:v}, v)$ 
7     else
8        $(t, \text{lev}([\mathbf{b}])) \leftarrow (\text{lev}([\mathbf{b}]), v)$ 
9        $\text{omg}([\mathbf{b}]) \leftarrow \text{omg}([\mathbf{b}]) \prod_{t < u \leq v} \varpi_{[\mathbf{b}]}^u$ 
10      while  $k < N$  do
11         $([\mathbf{x}_0], \pi_k, j) \leftarrow ([\mathbf{b}], 1, 0)$ 
12        while  $\rho([\mathbf{x}_j]) > r$  and  $[g_{1:v}]([\mathbf{x}_j]) \notin [\epsilon_{1:v}]$  do
13          if  $\left| \log_2 \left( \frac{\text{omg}([\mathbf{b}])}{\text{omg}([\mathbf{x}_j])} \pi_k \right) \right| > \sigma$  goto 31
14          ifundef  $\text{cuts}([\mathbf{x}_j]) \leftarrow \text{Cut}([\mathbf{x}_j], v)$ 
15           $([l_{j+1}], [r_{j+1}]) \leftarrow \text{cuts}([\mathbf{x}_j])$ 
16           $j \leftarrow j + 1$ 
17          ifundef  $(\text{omg}([r_j]), \text{lev}([r_j])) \leftarrow (\omega_{[r_j]}^{1:v}, v)$ 
18          else
19             $(t, \text{lev}([r_j])) \leftarrow (\text{lev}([r_j]), v)$ 
20             $\text{omg}([r_j]) \leftarrow \text{omg}([r_j]) \prod_{t < u \leq v} \varpi_{[r_j]}^u$ 
21          ifundef  $(\text{omg}([l_j]), \text{lev}([r_j])) \leftarrow (\omega_{[l_j]}^{1:v}, v)$ 
22          else
23             $(t, \text{lev}([l_j])) \leftarrow (\text{lev}([l_j]), v)$ 
24             $\text{omg}([l_j]) \leftarrow \text{omg}([l_j]) \prod_{t < u \leq v} \varpi_{[l_j]}^u$ 
25           $(\nu_{[l_j]}, \nu_{[r_j]}) \leftarrow (\text{omg}([l_j]), \text{omg}([r_j]))$ 
26           $[\mathbf{x}_j] \leftarrow \text{Bern}([l_j], \nu_{[l_j]}, [r_j], \nu_{[r_j]})$ 
27           $\pi_k \leftarrow \frac{\nu_{[\mathbf{x}_j]}}{\nu_{[l_j]} + \nu_{[r_j]}} \pi_k$ 
28        end
29         $w_k \leftarrow P(\mathbf{X} \in [\mathbf{x}_j]) / \pi_k$ 
30        Build  $\mathbf{y}_k$  by sampling  $[\mathbf{X} | \mathbf{X} \in [\mathbf{x}_j]]$ 
31         $k \leftarrow k + 1$ 
32        for  $i \leftarrow j$  to 1 do
33           $[\mathbf{x}_{i-1}] \leftarrow \text{Shrink}([l_i], [r_i])$ 
34           $\text{omg}([\mathbf{x}_{i-1}]) \leftarrow \text{omg}([l_i]) + \text{omg}([r_i])$ 
35        end
36      end
37 end

```

Algorithm 4: Incremental sampling.

Our purpose is to sample successively the sequence of conditional vectors:

$$[\mathbf{X} | g_{1:1}(\mathbf{X}) \in [\epsilon_{1:1}]], \dots, [\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]], \dots, [\mathbf{X} | g_{1:U}(\mathbf{X}) \in [\epsilon_{1:U}]], \quad (20)$$

in an efficient way. When the sampling is done by means of algorithm 3, the sampling structure produced for each conditional vector, $[\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]]$ where $v \in \llbracket 1, U \rrbracket$, is lost from a simulation process to another. This is inefficient, since it is foreseeable that the sampling structure for vector $[\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]]$ is certainly informative for sampling vector $[\mathbf{X} | g_{1:v+1}(\mathbf{X}) \in [\epsilon_{1:v+1}]]$.

Algorithm 4 is an update of algorithm 3, with the purpose of incrementally building the sampling structures for the sequence $([\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]])_{1 \leq v \leq U}$. Predictive weight is now incremental and takes multiplicative form:

$$\omega_{[\mathbf{x}]}^{1:v} = P(\mathbf{X} \in [\mathbf{x}]) \prod_{u=1}^v \varpi_{[\mathbf{x}]}^u. \quad (21)$$

Prefix $P(\mathbf{X} \in [\mathbf{x}])$ in (21) is not out of place: recall that $\omega_{[\mathbf{x}]}^{1:v}$ should approximate $P(\mathbf{X} \in [\mathbf{x}] \& g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]) = P(\mathbf{X} \in [\mathbf{x}]) P(g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}] | \mathbf{X} \in [\mathbf{x}])$. The multiplicative form is then a consequence of the incremental constraint.

In order to deal with the incremental construction of the samplers, algorithm 4 introduces new dictionary, `lev`, complementary to `cuts` and `omg`. Constraint level $t = \text{lev}([\mathbf{x}])$ indicates that weight `omg([\mathbf{x}])` has been actually computed with constraints sequence $g_{1:t}(\mathbf{X}) \in [\epsilon_{1:t}]$. Indeed, the sampling structure is not guaranteed to be up-to-date for all subtrees of the sampling structure, and it is thus necessary to trace this upgrade level. On the other hand, function call `Cut(v, [\mathbf{x}])`, which creates a cut of box $[\mathbf{x}]$, now depends on constraint level v . So, news cuts are created by considering all constraints information, although this will of course not be the case for cuts related to past levels $u < v$. At last, new function `Shrink` is introduced. Code $[\mathbf{x}] \leftarrow \text{Shrink}([\mathbf{l}], [\mathbf{r}])$ replaces $[\mathbf{x}]$ by the smallest box containing both $[\mathbf{l}]$ and $[\mathbf{r}]$; *this replacement is deep* and thus concerns dictionaries `cuts`, `omg` and `lev` as well. Even root variable $[\mathbf{b}]$ in algorithm 4 can have its value replaced. Nevertheless, deep replacements are cost-free here thanks to the way the structures are implemented.

In algorithm 4, parts taken unchanged from algorithm 3 are in black while evolutions are in blue. A `for` loop over the number of constraints, v , is added at lines 3 and 36 and repeats the entire sampling process for a given set of constraints. The sampling process itself does not change much, but incremental computations are now done on `omg` and `lev` is updated in consequence: lines 5 to 8, lines 16 to 19 and lines 20 to 23. Incremental computation means that undefined value `omg([\mathbf{x}])` is set equal to $\omega_{[\mathbf{x}]}^{1:v}$, while already defined value `omg([\mathbf{x}])` is multiplicatively completed by $\prod_{u=t+1}^v \varpi_{[\mathbf{x}]}^u$. At line 32, the shape of box $[\mathbf{x}_{i-1}]$ is recomputed from its cut. This is better for improving future incremental computation of `omg` during the sampling process.

3.3 Practical implementation and parallelism

Practical implementation. Algorithm 3 and 4 draw the main principles of our sampling methods. Some implementation details are described now.

Initial predictive weight. We implement the following definition of ω :

$$\omega_{[\mathbf{x}]} = \frac{\mu([g]([\mathbf{x}]) \cap [\boldsymbol{\epsilon}])}{\mu([g]([\mathbf{x}])))} P(\mathbf{X} \in [\mathbf{x}]) , \quad (22)$$

where μ is Borel measure on \mathbb{R}^m . This definition checks properties (15) and (16). It also tries a rough approximation for $P(\mathbf{X} \in [\mathbf{x}] \& g(\mathbf{X}) \in [\boldsymbol{\epsilon}])$. For incremental implementation, definition extends naturally as a product:

$$\omega_{[\mathbf{x}]}^{1:v} = P(\mathbf{X} \in [\mathbf{x}]) \prod_{u=1}^v \varpi_{[\mathbf{x}]}^u \text{ with } \varpi_{[\mathbf{x}]}^u = \frac{\mu([g_u]([\mathbf{x}]) \cap [\boldsymbol{\epsilon}_u])}{\mu([g_u]([\mathbf{x}])))} . \quad (23)$$

Cutting strategy. The definition of `Cut` is an important choice. Our algorithm selects a cut $([\mathbf{y}], [\mathbf{z}])$ randomly in regards to the following criteria:

- Favor cuts such that $\omega_{[\mathbf{y}]} \ll \omega_{[\mathbf{z}]}$ or $\omega_{[\mathbf{y}]} \gg \omega_{[\mathbf{z}]}$,
- Avoid overly elongated $[\mathbf{x}]$, i.e. such that $\frac{\max_i(x_i^+ - x_i^-)}{\min_i(x_i^+ - x_i^-)} \gg 1$,

In addition, some cuts history simplifications are implemented (next point).

Reducing sampling structure. Our implementation tries to optimize the structure of cut, weight, and level history, when it is possible. Different cases are implying such reductions, which are not mentioned here. A typical example is shown in figure 3. Assume that $([\mathbf{y}], [\mathbf{z}])$ is a cut of $[\mathbf{x}]$, $([\mathbf{t}], [\mathbf{u}])$ is a cut of $[\mathbf{z}]$ and $[g]([\mathbf{u}]) \cap [\boldsymbol{\epsilon}] = \emptyset$, i.e. $[\mathbf{u}]$ is outside the subpaving and its border. Then, boxes $[\mathbf{z}]$ and $[\mathbf{u}]$ are useless and should be removed from the structure.

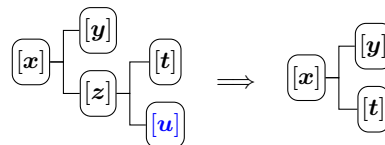


Fig. 3: Structure reduction

Discarding first samples. M first samples are discarded, so as to initialize the structure of the sampler. After that, N samples are sampled and returned.

Setup. The algorithms are rather simple to set up. Except for the choice of ω , which is structural, r , M and σ are the only parameters to be defined.

Parallelization. This paragraph gives additional information on how incremental algorithm 4 is implemented with parallelization. Data structures `cuts`, `omg` and `lev` deal rather well with parallel processing, so that our implementation is multithread. Nevertheless, the practical implementation needs some work.

In our implementation, the dictionaries `cuts`, `omg` and `lev` are stored in a same structure. This structure is composed of shared references to vertices which are stored within a Slab allocation [2]. Slab allocation is a structure, which allows easy and efficient deallocation and reallocation, while reducing the risk of fragmentation of the memory. A vertex is an element of the history of the sampling process, which contains the following fields:

- A box $[\mathbf{x}]$ with its weight $\text{omg}([\mathbf{x}])$ and its level $v = \text{lev}([\mathbf{x}])$, which is an indicator that $\text{omg}([\mathbf{x}])$ has been computed from $g_{1:v}$ and $[\epsilon_{1:v}]$,
- The type of vertex: a node or a leaf? Is property $[g_{1:v}](\mathbf{x}) \subset [\epsilon_{1:v}]$ true?
- When vertex is a node and box $[\mathbf{x}]$ has a cut $([\mathbf{l}], [\mathbf{r}]) = \text{cuts}([\mathbf{x}])$, then the vertex is completed with weak references to the vertices related to boxes $[\mathbf{l}]$ and $[\mathbf{r}]$ respectively.

Shared references allow multiple threads to access a same data; moreover, the existence of a shared reference will prevent the erasure of referenced data. Weak references are derived from shared references and also allow threads to access the data, but they do not prevent the erasure of referenced data. By restricting the shared references to the Slab allocation, we prevent the possibility of memory leak, especially when the sampling structure is reduced by a thread. Since the structure is shared among threads, the access to the vertices are protected by memory locks: the locks allow either many read-only accesses or only one read-write access to the data.

This data structure allowed some fair performance increase in comparison to conference paper [4]. The memory management is now smoother and more efficient. In particular, we no longer have process freezing phenomena following heavy memory management.

The definition of an efficient data structure is a main achievement when designing a multithread implementation. Now, we give some clarification on how the sampling process is shared between the threads:

- Multiple instances (usually as many as the number of processor threads) of the sampling loop are run simultaneously. This includes descending through the sampling structure, possible extension of this structure, possible sampling of box and point if conditions are met, and backward reconstruction of weights and restructuring of the sampling structure,
- If a vertex currently processed by an instanced loop has been deleted by another instance, then the vertex subtree is first deleted by the process, and afterward, the loop is restarted.
- Processes are stop when the desired number of samples is obtained,

4 Application to Bayesian optimization

Bayesian optimization is a key technique for black-box optimization, and it is actually a great motivation for this work.

Assume that one needs to optimize a function which is not well known, and which may be computed by a highly costly process (a heavy simulation, tests made by human teams on the grounds, etc.). Of course the optimization should be made by sparing at best the number of calls to the costly evaluation.

In [8], Jones, Schonlau and Welch proposed the efficient global optimization method (EGO) for addressing such kind of problem. The idea is to use a surrogate model under the form of a functional random variable. This functional random variable is described by means of a Gaussian process with correlation depending

on spatial distance (kriging). Based on such modelling, the construction of an *optimal* parameter sequence to be evaluated is obtained by iterating:

- Compute the posterior law of the functional random variable, according to the past evaluations,
- Compute the expected improvement function, in regards to the posterior law and the already best computed value. This expected improvement is an indicator of the next parameter to be evaluated,
- Find parameter optimizing the expected improvement and evaluate it.

The approach relies on the construction of the Gaussian modelling, and it tends to be less efficient when the dimension of the optimization space increases. Works have been made in order to deal with this dimensional issue and improve the modelling; *e.g.* [6]. But this issue is still a true challenge.

Beside, EGO itself is a form of Bayesian optimization introduced by Mockus [10, 9, 13]. We proposed in [3] a Bayesian optimization approach in the context of a nonlinear function depending on a model noise; the realization of the noise was the unknown information of the problem. We hoped by a nonlinear approach to better handle some of the difficulty induced by the dimension. The work was not complete, since we were not able to build a good conditional sampling at that time. The purpose of this paper and of [4] was to build such sampler.

It is beyond the scope of this chapter to detail the seminal works of Mockus, or those of Jones, Schonlau and Welch. We will introduce the subject by a short description of EGO algorithm. Then, we present our nonlinear approach.

Formalization: Function $\gamma \mapsto g(\gamma, \mathbf{x}^o)$ is to be minimized. Parameter \mathbf{x}^o is unknown, but \mathbf{x}^o is a realization of random vector \mathbf{X} , whose law $F_{\mathbf{X}}$ is known. In order to optimize γ , we are allowed to request an evaluator for computing $g(\gamma, \mathbf{x}^o)$, but each call to this evaluator is costly. The objective is then to solve $\gamma^o \in \arg \min_{\gamma} g(\gamma, \mathbf{x}^o)$ by optimizing at the same time the sequence of evaluated cases γ_u and of their evaluations $e_u = g(\gamma_u, \mathbf{x}^o)$.

EGO method: Main idea of EGO consists in approximating $g(\gamma, \mathbf{X})$ by a Gaussian process (GP) $\gamma \mapsto G(\gamma)$. Given past evaluated cases γ_u and evaluations $e_u = g(\gamma_u, \mathbf{x}^o)$ for $u \leq v$, posterior variable $[G(\gamma) | \forall u \leq v, G(\gamma_u) = e_u]$ is Gaussian and mathematically computed. This conditional variable is a main ingredient for computing the expected improvement, $EI(\gamma)$, which is an indicator of value γ for which next evaluation of $g(\gamma, \mathbf{x}^o)$ is promising:

$$EI(\gamma) = E_{G(\gamma) | \forall u \leq v, G(\gamma_u) = e_u} \max\{m_{\gamma} - G(\gamma), 0\} \text{ where } m_{\gamma} = \min_{1 \leq u \leq v} e_u. \quad (24)$$

It is computed mathematically. Since $EI(\gamma)$ indicates where evaluations are promising, it is maximized in order to chose next case $\gamma_{v+1} \in \arg \max_{\gamma} EI(\gamma)$ and request its evaluation $e_{v+1} = g(\gamma_{v+1}, \mathbf{x}^o)$.

Nonlinear algorithm: We took inspiration of EGO, but now, random function $g(\cdot, \mathbf{X})$ is used directly instead of GP approximation. Parameters γ_v and evaluations $e_v = g(\gamma_v, \mathbf{x}^o)$ are optimized by iterating algorithm 5. While $EI(\gamma)$ is formed differently in this algorithm, it is essentially the same as in definition (24) modulo an affine transform. It is theoretically defined by:

$$EI(\gamma) = E_{\mathbf{X} | \forall u \leq v, g(\gamma_u, \mathbf{X}) = e_u} \min\{g(\gamma, \mathbf{X}), m_\gamma\} \quad \text{with} \quad m_\gamma = \min_{1 \leq u \leq v} e_u. \quad (25)$$

In practice, $EI(\gamma)$ is obtained by a Monte Carlo derived from an approximated sampling of $[\mathbf{X} | \forall u \leq v, g(\gamma_u, \mathbf{X}) = e_u]$. *As a crucial ingredient:*

- Vector $[\mathbf{X} | \forall u \leq v, g(\gamma_u, \mathbf{X}) = e_u]$ is approximated by $[\mathbf{X} | g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]]$, where $g_u = g$ and ϵ_u is a small box around e_u .
- The sampling is done by the way of this approximation.

This approximated random vector, conditional to a subvariety, incorporates incrementally all constraints related to the past evaluations. The sampling is done by algorithm 3 or preferably by incremental algorithm 4.

Algorithm thus consists of:

1. Generating samples of posterior random vector $[\mathbf{X} g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]]$,	1	Function	<i>Process next measure</i>	
2. Building a Monte Carlo approximation of $EI(\gamma)$,	2	input	γ_u and $e_u \triangleq g(\gamma_u, \mathbf{x}^o)$ for $1 \leq u \leq v$	
3. Choosing parameter γ_{v+1} by minimizing $EI(\gamma)$ and evaluating it.	3	output:	γ_{v+1} and $e_{v+1} \triangleq g(\gamma_{v+1}, \mathbf{x}^o)$	
	4		Make samples of $[\mathbf{X} g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]]$	
	5		Compute $m_\gamma = \min_{1 \leq u \leq v} e_u$ and:	
	6		$EI(\gamma) \simeq E_{\mathbf{X} g_{1:v}(\mathbf{X}) \in [\epsilon_{1:v}]} \min\{g(\gamma, \mathbf{X}), m_\gamma\}$	
			Compute $\gamma_{v+1} \in \arg \min_\gamma EI(\gamma)$	
			Compute $e_{v+1} \triangleq g(\gamma_{v+1}, \mathbf{x}^o)$	
		end		

Algorithm 5: Sampler-based Bayesian optimizer

Implementations of step 4 are not detailed here. Essentially, we used a meta-heuristic method, the cross-entropy algorithm for optimization [12], for this task; this method is related to parameterized rare-event simulation.

5 Examples and tests

Regarding the conference publication [4], all the results presented here have been recomputed or supplemented by recent implementations of the algorithms.

5.1 Simulation: test cases

The tests presented here are performed for sampling algorithm 3 or its incremental version 4. The algorithms have been implemented in `Rust` language (www.rust-lang.org) and were processed on 7 threads. The algorithms are tested on mathematically simple simulation problems, in order to make the statistics of the results clear enough to analyze.

Parameters. All simulations have been achieved with the following parameters:

- $r = 0.001$ is the radius bound for second **while** stop condition,
- $M = 5000$ is the number of samples discarded during sampler initialization,
- $N = 50000$ is the number of sampled particles,
- $\sigma = 10$ on all tests.

Test Cases. Thorough the section, it is assumed that \mathbf{X} follows the uniform law on $\mathbf{b} = [-2, 2]^n$ with $n \in \{2, 3, \dots, 11\}$. Three cases are investigated:

Case (a): Are defined $g_a(\mathbf{x}) = \|\mathbf{x}\|_2 = \sqrt{\sum_{j=1}^n x_j^2}$ and $[\epsilon_a] = [0.95, 1.05]$. Then $g_a^{-1}([\epsilon_a])$ is a hyper-spherical shell, which approximates the unit hypersphere of dimension $n - 1$.

Case (b): For $n = 11$ and $0 \leq v \leq 9$, it is considered:

$$g_b(\mathbf{x}) = \left[\|\mathbf{x}\|_2 \quad x_3 \quad \dots \quad x_{2+v} \right] \text{ and } [\epsilon_b] = [\epsilon_a] \times [\mathbf{z}]^v, \quad (26)$$

with $[\mathbf{z}] = [-0.05, 0.05]$. Similarly to **(a)**, but with additional constraints, $g_b^{-1}([\epsilon_b])$ approximates a hypersphere of dimension $n - 1 - v$. When $v = 0$, we are back to case **(a)** with $n = 11$. When $v = 9$, then $g_b^{-1}([\epsilon_b])$ approximates the unit circle \mathcal{C} within the first two dimensions, related to coordinates x_1, x_2 .

Case (c): For $n = 11$ and $0 \leq v \leq 9$, it is considered:

$$g_c(\mathbf{x}) = \left[\|\mathbf{x}\|_2 \quad \min(|x_1|, |x_2|) \quad x_3 \quad \dots \quad x_{2+v} \right], \quad (27)$$

and $[\epsilon_c] = [\epsilon_a] \times [\alpha] \times [\mathbf{z}]^v$ with $[\alpha] = [0, 0.5]$. Thus, this case is obtained by adding constraint $0 \leq \min(|x_1|, |x_2|) \leq 0.5$ to subcases of **(b)**. For subcase $v = 9$ especially, we are approximately sampling on the (disjoint) union of the 4 subsegments $\mathcal{A}_1, \dots, \mathcal{A}_4$ of the unit circle \mathcal{C} , defined by:

$$\mathcal{A}_j = \left\{ [x_1, x_2] \in \mathcal{C} \middle/ \arg([x_1, x_2]) \in j \frac{\pi}{2} + \left[-\frac{\pi}{6}, \frac{\pi}{6}\right] \bmod 2\pi \right\}. \quad (28)$$

These 4 subsegments have the same size so that their probabilities are the same in regard to X .

Purpose of the Test Cases. Subsequently, case **(a)** is used in order to evaluate the performance of the sampling process both in accuracy and in efficiency for different dimensions. Case **(b)** is used in order to evaluate the efficiency of the sampling when the number of constraints increases; algorithm 3 and incremental algorithm 4 are also compared. Case **(c)** is used in order to evaluate the accuracy of the sampling in case of complex constraints which introduce disjoint modes.

5.2 Simulation: case (a)

This case is mathematically easy to predict. In [5], Dezert and Musso proposed a method, which may be used for uniformly sampling on an ellipsoid shell. Whatever, one must keep in mind that our approach is generic and can be applied to an infinite number of configurations.

Histograms. Figure 4 presents some histograms for case (a). For each subcase $n \in \{3, 7, 11\}$, we have computed the radius of all samples \mathbf{x} and built the associated histograms – subfigures (1), (3) and (5). For each $n \in \{3, 7, 11\}$ and for all $1 \leq i < j \leq n$, we have computed the angle of all samples (x_i, x_j) and built the associated histograms. From these $\frac{n(n-1)}{2}$ histograms of each subcase, we have computed the minimal, mean and maximal histograms. The results are shown in blue, green and red, respectively, and provide an hint on the error of the estimation – subfigures (2), (4) and (6).

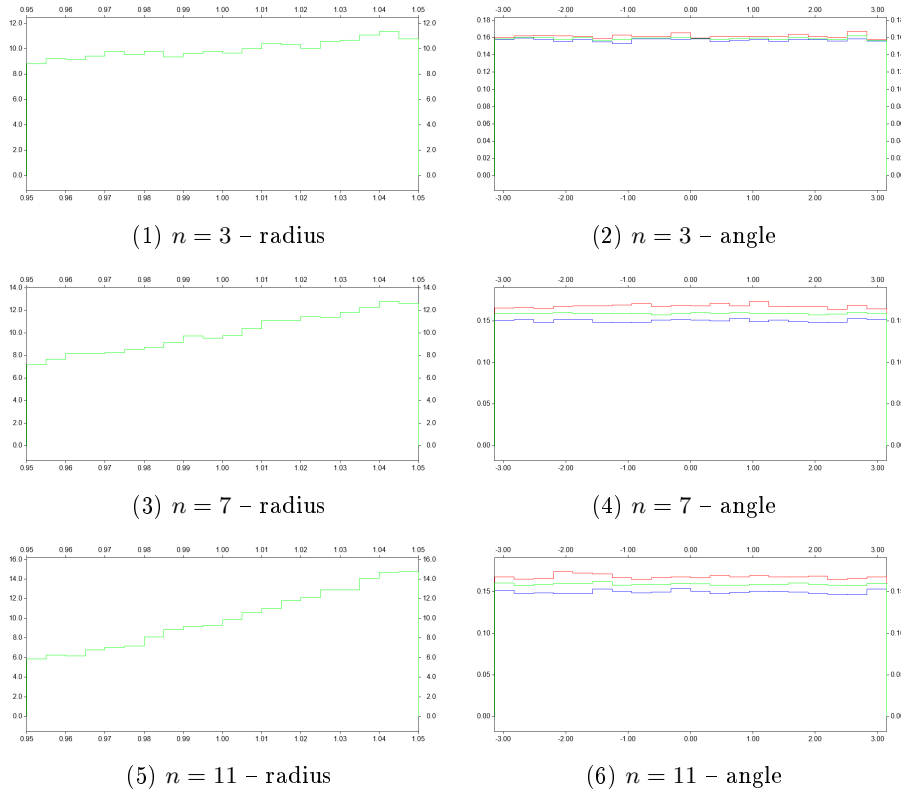


Fig. 4: Case (a) – histograms – 20 divisions

By symmetry, the *theoretical angle histograms* are uniform. The errors should be of the order of $\sqrt{\frac{20}{50000}} = 0,02$. In comparison, the errors figured in the angular histograms are quite acceptable, even for the highest dimension. The most interesting point is that there is no rupture in the histogram, which shows that the sampler does manage the subvariety structure. We do not have an error estimation for the radius histograms. Actually, the local probability should theoretically increase with the radius, this property being accentuated with the

dimension. This is what is obtained on the histograms. It is noteworthy however that the sides of these histograms are subject to additional errors implied by the border of the subvariety.

Process Statistics: Figure 5 presents some statistics graphs for case (a).

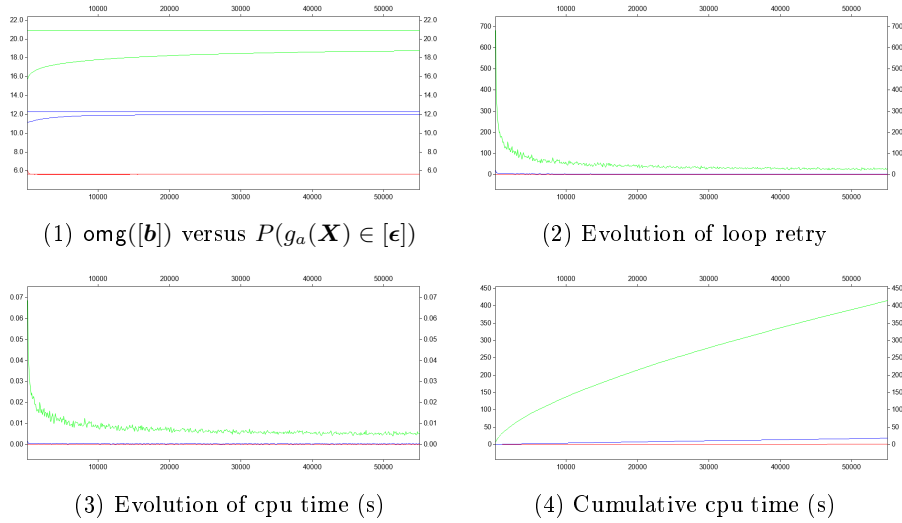


Fig. 5: Case (a) – statistics – $n = 3, 7, 11$

Results are plotted against the number of generated samples $k = 1 : 55000$. Cases $n = 3, 7, 11$ are drawn with respective colors, red, blue and green.

Subfigure (1) shows how value $-\log_2(\text{omg}([\mathbf{b}]))$ evolves and approximates $-\log_2(P(g_a(\mathbf{X}) \in [\epsilon]))$. Each curve increases to theoretical value (same color line), but performance decreases with dimension.

Subfigures (3) and (4) present the cumulative cpu-time per thread and the evolution of the cpu-time per sample and per thread consumed by the process (expressed in second). In comparison to [4], there is no more time discontinuity caused by intermittent memory allocations. Our memory management has been improved as shown in section 3.3. We notice clearly that the sampling efficiency increases with the number of generated samples. However, the cumulative cpu time still increases dramatically with the dimension (the memory use evolves similarly). Although the curse of dimension has been delayed by our approach, it is still there.

The number of loop retries during the sampling is plotted in subfigure (2). It is an interesting indication of the achievement of the sampling structure. It decreases with the number of samples and becomes small, even for the highest dimension (around 28 for $n = 11$). This result should be compared to the probability of the subvariety (around 10^{-6} for $n = 11$).

5.3 Simulation: case (b)

Figure 6 presents synthetic performance curves, as well as comparative histogram examples for algorithm 3 and algorithm 4 on case (b).

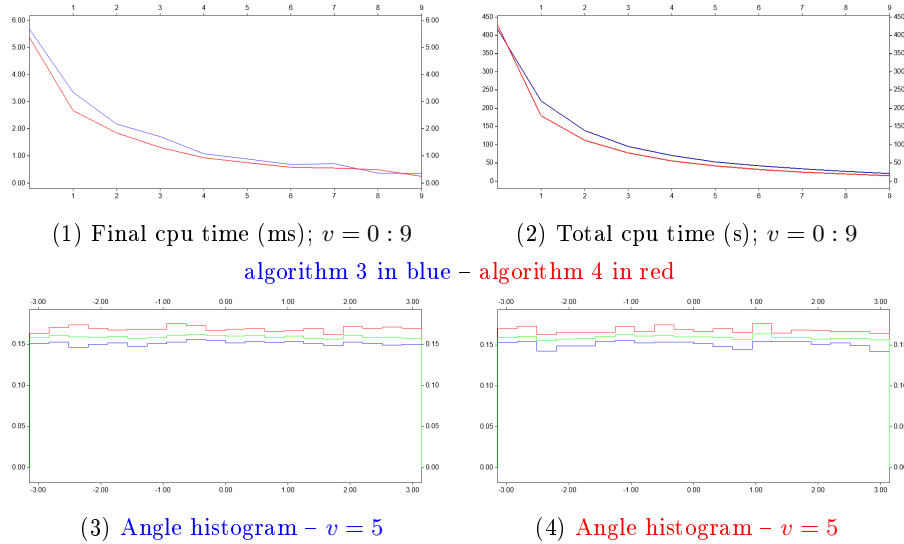


Fig. 6: Case (b) – speed and histograms

Subfigures (1) and (2) are plotted against the number of additional constraints $v = 0 : 9$. Subfigure (1) presents final cpu-time (in ms) per sample and per thread consumed by algorithm 3 (in blue) and algorithm 4 (in red). Both curves decrease and approach to zero. Subfigure (2) presents total cpu-time (in s) per thread consumed by these algorithms. Both curves decrease. Red curves are dominated by blue curves, which means better efficiency of algorithm 4 against algorithm 3.

Subfigure (3) presents angle histogram for subcase $v = 5$ from computations by algorithm 3. Subfigure (4) presents angle histogram for same subcase from computations by algorithm 4. Both results are comparable in terms of accuracy. The theoretical histogram is uniform, which conforms to these results.

As a conclusion here, the performance of the sampler is likely to increase with the number of constraints, and this is a useful quality. Incremental algorithm 4 is more efficient than algorithm 3, while ensuring equivalent accuracy. Algorithm 4 must therefore be favored for incremental constraints. In last example of section 5.5 on Bayesian optimisation, implementation of algorithm 4 was even necessary for an acceptable computation time.

5.4 Simulation: case (c)

Figure 7 presents computed radial histograms and angular histograms (for coordinates x_1 and x_2 – there is only one) for subcases $v = 6$ and $v = 9$.

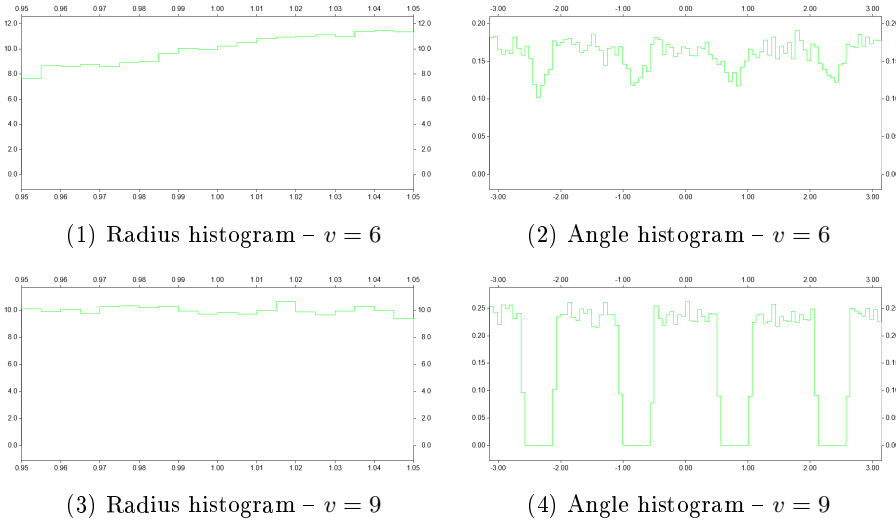


Fig. 7: Case (c) – histograms – radius: 20 divisions; angle: 100 divisions

Subfigure (3) presents radial histogram for final case $v = 9$. Subfigure (4) presents angle histogram for final case $v = 9$. The quality of the histograms is comparable to what was observed previously. Due to the constraint configuration of subcase $v = 9$, the *theoretical radius histogram* is uniform, and the *theoretical angular histogram* is uniform around each subsegment, \mathcal{A}_1 to \mathcal{A}_4 , with a gradual decrease on the borders. The generated histograms actually comply with these properties.

Subfigures (1) and (2) present intermediate subcase $v = 6$. Since some other coordinates than x_1 and x_2 are relaxed, constraint $x_1^2 + x_2^2 \simeq 1$ is also relaxed. Actually, theoretical angular histogram is even no longer discontinuous around subsegments, \mathcal{A}_i , for this reason.

As a preliminary conclusion, we consider that our sampling method is globally performant in sampling conditionally to subvarieties. A future issue will be to increase the dimension of the sampling space.

5.5 Bayesian optimization

A simple geometric problem. We intend to find the isobarycenter $\gamma = (a, b)$ of 4 *unknown points* $M_i = (x_{2i-1}^o, x_{2i}^o) \in [-5, 5]^2$ with $i \in \{1, 4\}$. The only

approach that is possible for us is to test some solutions by requesting for a costly measurement; for requirement $\gamma_u = (a_u, b_u)$, this measurement evaluates:

$$e_u = g(\gamma_u, \mathbf{x}^o) = \|\gamma_u - h(\mathbf{x}^o)\|_2 \quad \text{with} \quad h(\mathbf{x}^o) = \frac{1}{4} \sum_{i=1}^4 M_i = \frac{1}{4} \sum_{i=1}^4 (x_{2i-1}^o, x_{2i}^o). \quad (29)$$

Our purpose is to optimize (a, b) by minimally requesting evaluation $g(\gamma, \mathbf{x}^o)$.

Geometric solution. Each measure restricts the solution to a circle. After 2 measures, we usually have to choose between two points, and the solution is found equiprobably at step 3 or 4.

Tests and Results. Points M_1, \dots, M_4 are $(2, -1)$, $(3, 2)$, $(-\frac{3}{2}, 4)$, $(\frac{1}{2}, 3)$. Their isobarycenter is $(1, 2)$. We used a sampler with $M = 5000$, $N = 10000$ and $[\epsilon] = [-\frac{1}{100}, \frac{1}{100}]^u$. Variable \mathbf{X} is considered uniform on $[-5, 5]^8$. Table 1 presents a typical optimized sequence for parameters $\gamma_u = (a_u, b_u)$. Optimization of EI is done the cross-entropy method [12] and is thus near-optimal. In this example,

Table 1: Exemple of optimization sequence

u	1	2	3	4	5	6	7	8	9	10	11	12	13
a_u	0.01	-0.41	0.84	0.97	0.99	1.00	0.98	1.02	-3.79	0.98	0.98	-0.39	1.03
b_u	0.02	-1.94	2.05	1.99	2.02	1.98	2.02	1.99	6.49	2.00	1.99	5.40	2.00
e_u	2.21	4.19	0.16	0.03	0.02	0.02	0.03	0.02	6.57	0.02	0.03	3.67	0.03

best value for e_u is 0.02. Best values are generally found around 0.02: this is a consequence of error interval $[\epsilon]$, which is not zero size.

Table 2: Results for 100 runs

u_o	1	2	3	4	5	6	7	8	9	10	11	12	13
$e_{u_o} \leq 0.02$	0%	0%	26%	40%	10%	9%	2%	3%	3%	3%	1%	1%	1%
$e_{u_o} \leq 0.04$	0%	0%	36%	57%	6%	1%	0%	0%	0%	0%	0%	0%	0%
$e_{u_o} \leq 0.09$	0%	0%	40%	57%	3%	0%	0%	0%	0%	0%	0%	0%	0%

Table 2 presents the results for 100 runs. For these runs, the mean cpu time is 2627s, where *non-incremental* algorithm 3 has been used for simulation. Three convergence conditions, $e_u \leq e_{cond}$ with $e_{cond} \in \{0.02, 0.04, 0.09\}$, are considered for these runs. Value $u_o = \min\{u/e_u \leq e_{cond}\}$ is the final step for which convergence condition is met. Percentage is given for each final step value u_o . In all cases, we notice that most optimal values are found at steps 3 or 4. The result tends to be equilibrated on 3 and 4 when stopping criterion is relaxed.

EGO method. For the sake of comparison, we apply EGO to our problem. Gaussian process G with covariance $\text{cov}(G(\gamma), G(\gamma')) = \frac{1}{4} \exp(-\frac{1}{2} \|\gamma - \gamma'\|_2)$ and mean 0 is used. Table 3 summarizes the results of 100 tests. Each test imple-

Table 3: Results for EGO method

e_{cond}	0.5	0.2	0.1	0.05	0.02	0.01	0.005	0.002	0.001
% : $u_o \leq 100$	91%	74%	69%	64%	55%	46%	42%	29%	14%
worse(u_o)	55	78	96	97	97	84	93	93	99
mean(u_o)	11.8	19.4	23.1	25.5	30	31.3	37.5	56.3	60.5
best(u_o)	1	1	8	8	8	15	15	19	19

ments 100 successive evaluations. Convergence conditions are $e_u \leq e_{cond}$ with different values for e_{cond} ranging from 0.5 to 0.001. Again, value u_o is the final step for which convergence condition is met. For each value e_{cond} , the table indicates the percentage of tests which succeeded to reach convergence with less than 100 evaluations. In case of convergence, the worse, mean and best final steps are given. EGO is outperformed here. However, EGO is able to achieve more refined results – e.g. $e_{cond} = 0.001$ – if the evaluation budget is relaxed: indeed, EGO is based on exact mathematical computation of the expected improvement.

Decision space of higher dimension. Variable \mathbf{X} is uniform on $[-5, 5]^8$ with realization $\mathbf{x}^o = (2, -1, 3, 2, -\frac{3}{2}, 4, \frac{1}{2}, 3)$. Points $M_{i,j} = (m_{i,j,k})_{1 \leq k \leq 28}$ with $1 \leq j < i \leq 8$ are built from \mathbf{X} by setting $m_{i,j,k} = X_i X_j$ if $k = j + \frac{(i-2)(i-1)}{2}$ and $m_{i,j,k} = 0$ else. Function $g(\gamma, \mathbf{X})$ is defined by $g(\gamma, \mathbf{X}) = \|\gamma - h(\mathbf{X})\|_2$ where $h(\mathbf{X}) = \frac{1}{28} \sum_{i=2}^8 \sum_{j=1}^{i-1} M_{i,j}$. The problem is thus of much higher dimension.

We had some theoretical issues on this example. The Bayesian estimation was unable to resolve some remaining punctual cases for zero dimension subvarieties. This difficulty has been overcome by adding exclusion constraints around already tested cases. Table 4 gives an example of run (progress is in bold). Incremental

Table 4: Example in a space of higher dimension

u	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
e_u	1.02	1.04	1.18	0.77	0.68	0.68	0.61	1	0.33	0.12	0.18	0.21	0.09	0.09	0.09	0.09	0.05	0.10	0.09	0.12

algorithm 4 were needed for simulation, and optimization has been performed within a day. On the other hand, algorithm 3 was far too slow in this case for optimizing in reasonable time.

6 Conclusion

We proposed an original dichotomous method for sampling a random vector conditionally to a subvariety. This method has been parallelized, and we proposed

non-incremental and incremental implementations. This generic approach, inspired from interval analysis, is accurate and efficient up to a space of dimension 11. We have shown how it could be applied efficiently to Bayesian optimization problems. The work is promising from theoretical and applicative point of view and offers some improvement perspectives. A main issue is to enhance the efficiency of the approach with respect to higher dimensions. With this perspective in mind, some relaxation techniques applied to the subvariety may be considered.

References

1. Alefeld, G., Mayer, G.: Interval analysis: theory and applications. *Journal of Computational and Applied Mathematics* **121**(1), 421–464 (2000)
2. Bonwick, J., Microsystems, S.: The slab allocator: An object-caching kernel memory allocator. In: *USENIX Summer*. pp. 87–98 (1994)
3. Dambreville, F.: Optimizing a sensor deployment with network constraints computable by costly requests. In: Thi, H.A.L., Dinh, T.P., Nguyen, N.T. (eds.) *Modelling, Computation and Optimization in Information Systems and Management Sciences. Advances in Intelligent Systems and Computing*, vol. 360, pp. 247–259. Springer (2015)
4. Dambreville, F.: Simulating a Random Vector Conditionally to a Subvariety: A Generic Dichotomous Approach. In: *SIMULTECH'21, 11th International Conference on Simulation and Modeling Methodologies, Technologies and Applications*. pp. 109–120 (Jul 2021)
5. Dezert, J., Musso, C.: An efficient method for generating points uniformly distributed in hyperellipsoids. In: *Workshop on Estimation, Tracking and Fusion: A Tribute to Bar-Shalom*. Monterey, California (May 2001)
6. Hebbal, A., Brevault, L., Balesdent, M., Taibi, E.G., Melab, N.: Efficient global optimization using deep gaussian processes. In: *2018 IEEE Congress on Evolutionary Computation (CEC)*. pp. 1–8 (2018)
7. Jaulin, L., Kieffer, M., Didrit, O., Walter, E.: *Applied Interval Analysis with Examples in Parameter and State Estimation, Robust Control and Robotics*. Springer London Ltd (Aug 2001)
8. Jones, D., Schonlau, M., Welch, W.: Efficient global optimization of expensive black-box functions. *Journal of Global Optimization* **13**, 455–492 (12 1998). <https://doi.org/10.1023/A:1008306431147>
9. Mockus, A., Mockus, J., Mockus, L.: Bayesian approach adapting stochastic and heuristic methods of global and discrete optimization. *Informatica (lithuanian Academy of Sciences)* **5**, 123–166 (1994)
10. Močkus, J.: On bayesian methods for seeking the extremum. In: Marchuk, G.I. (ed.) *Optimization Techniques IFIP Technical Conference Novosibirsk, July 1–7, 1974*. pp. 400–404. Springer Berlin Heidelberg, Berlin, Heidelberg (1975)
11. Morio, J., Balesdent, M., Jacquemart, D., Vergé, C.: A survey of rare event simulation methods for static input–output models. *Simulation Modelling Practice and Theory* **49**, 287–304 (2014)
12. Rubinstein, R.Y., Kroese, D.P.: *The Cross Entropy Method: A Unified Approach To Combinatorial Optimization, Monte-Carlo Simulation (Information Science and Statistics)*. Springer-Verlag (2004)
13. Zilinskas, A., Zhigljavsky, A.: Stochastic global optimization: A review on the occasion of 25 years of informatica. *Informatica* **27**, 229–256 (01 2016)