

SMGO: A Set Membership Approach to Data-Driven Global Optimization

Lorenzo Sabug Jr.*, Fredy Ruiz, Lorenzo Fagiano**

*Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano
Piazza Leonardo da Vinci 32, 20133 Milano, Italy*

Abstract

Many science and engineering applications feature non-convex optimization problems where the objective function can not be handled analytically, i.e. it is a black box. Examples include design optimization via experiments, or via costly finite elements simulations. To solve these problems, global optimization routines are used. These iterative techniques must trade-off exploitation close to the current best point with exploration of unseen regions of the search space. In this respect, a new global optimization strategy based on a Set Membership (SM) framework is proposed. Assuming Lipschitz continuity of the cost function, the approach employs SM concepts to decide whether to switch from an exploitation mode to an exploration one, and vice-versa. The resulting algorithm, named SMGO (Set Membership Global Optimization) is presented. Theoretical properties regarding convergence and computational complexity are derived, and implementation aspects are discussed. Finally, the SMGO performance is evaluated on a set of benchmark non-convex problems and compared with those of other global optimization approaches.

1. Introduction

In many science and engineering fields, such as mechanical design, fluid-dynamics, electromagnetics, multi-physics simulations, control systems tuning, and chemical experiments, black-box optimization problems arise. Black-box functions are named as such due to the fact that an explicit mathematical model is unavailable, or too complicated to be handled analytically. In these cases, the optimization strategy can rely only on function values obtained through empirical tests. Moreover, these problems may present several local minima, and the time and resources required to carry out a single function evaluation are rather large, so that the solution method shall make the most efficient use of

the available trials. To highlight this fact, we refer to the process of performing a new test as *long function evaluation*.

Seeking a minimizer in this framework is referred to as black-box optimization [1], or derivative-free optimization [2]. This is in contrast for example to gradient-based and Newton-type methods where the (first, or even further) derivatives are also assumed to be available. If the black-box function is assumed differentiable, one could resort to gradient estimation techniques to still employ these methods, which however require a rather large number of long function evaluations to estimate a local quantity (i.e., the gradient at the currently evaluated point) [3] and are designed to converge to a local optimum by always improving from the initial guess, not to explore the decision space searching for a global solution.

Due to time/resource limitations, in global optimization there are two conflicting aspects that must be considered when choosing the next test point for a long function evaluation [3]. The first, *exploitation*, pertains to improving the result by testing more points in the vicinity of the current best known one. The second, *exploration*, aims to gain

*Corresponding author. The first author would like to acknowledge the support of the Department of Science and Technology–Science Education Institute (DOST-SEI) of the Philippines for his research.

**This research has been supported by the Italian Ministry of University and Research (MIUR) under the PRIN 2017 grant n. 201732RS94 "Systems of Tethered Multicopters".

Email addresses: lorenzojr.sabug@polimi.it (Lorenzo Sabug Jr.), fredy.ruiz@polimi.it (Fredy Ruiz), lorenzo.fagiano@polimi.it (Lorenzo Fagiano)

more information about the cost function in other regions of the search space, seeking other, possibly better, minima. Exploitation and exploration strategies are at the core of most global optimization approaches, including stochastic global search, response surface, and Lipschitz-based methods.

Stochastic global search algorithms, such as random search methods [4], particle swarm algorithms [5], and its offshoot variants [6, 7], use a population of search points per iteration. According to the corresponding function values, heuristics are applied to calculate the search point locations for the next iteration (or “generation”). These methods are widespread but require numerous long function evaluations per generation, which can be impractical. Moreover, aspects such as guaranteed convergence and optimality gap may be difficult to analyze. Response surface methods, on the other hand, generate an approximation of the black-box function via, e.g., Gaussian process regression [1], radial basis functions [8, 9], or neural networks [10]. These approximations can then be used for the exploitation or exploration routines, whichever is appropriate. Albeit the number of long function evaluation is lower in this case, optimality gap is again not easily assessed, and there is an additional exponential computational burden to derive and refine the approximated objective function.

Another category of black-box optimizers are so-called Lipschitz-based algorithms. They rely on the assumption that the long function is Lipschitz continuous. An early method, proposed independently by Piyavskii [11] and Shubert [12], uses a known Lipschitz constant to draw the lower bounds of the black-box function given the existing samples, starting with the bounds (corners) of the search space. In this context, the Piyavskii-Shubert (“sawtooth”) method chooses the next point by searching the location of the minimum lower bound. However, the method is proposed for one-dimensional function optimization only; furthermore, Lipschitz constants are rarely known in practical cases, and need to be estimated in the process.

A later proposal named DIRECT [13, 14, 15] is a modification to Piyavskii-Shubert method, removing the need for a known Lipschitz constant and starting samples from the corners. Instead, it invariably performs a first sample at the center of the search space, which is afterwards subdivided into three intervals. The Lipschitz constant estimate is used to draw the lower bounds, which are used to select potentially optimal intervals (hyperrectan-

gles). Such hyperrectangles are iteratively subdivided, and their corresponding centers are sampled by batch, and potential intervals for sampling are selected again. This cycle proceeds until the maximum number of iterations is reached, with each iteration performing multiple long function evaluations. An offshoot method named DISIMPL [16, 17] employs simplices instead of hyperrectangles, thus extending the usage from rectangular search spaces only to more general polytopic ones. On one hand, DIRECT and DISIMPL are completely deterministic and reproducible (assuming no noise on function evaluation), however, they only admit a fixed starting point, and do not propose how to handle existing sampled points at the beginning of the algorithms. Furthermore, the batch-based sampling means that the next set of trial points is not chosen until the current batch of sampling is finished, which can mean overlooking possible precious information added by individual samples as they come in.

LIPO and AdaLIPO [18] use a random binary variable to decide between exploitation and exploration; the only difference being that the former assumes a fixed Lipschitz constant, while the latter estimates it from existing data. For exploitation, they choose a random point within the set of potential optimizing points, i.e. those whose lower bound on the black-box function is below the current best sample. For exploration, it chooses a random point from the entire search space. These algorithms, owing to the simple implementation, are fast in terms of choosing the next sampling point, and can easily scale to high dimensional problems without significant hit on computational burden. However, such a completely randomized exploitation/exploration sampling may be inefficient in finding new local minima, because it does not fully utilize the shape of lower bounds. Furthermore, the results of such an algorithm are unrepeatable because of its randomized design.

In the described context, we propose a new Lipschitz-based algorithm that exploits Set Membership (SM) nonlinear function approximation theory [19]. So far, SM theory proved effective for nonlinear system identification [20], filter design [21], and controller design [22], among others. Here, we employ SM theory for the first time in global optimization. Our paper delivers the following contributions:

- The Set Membership Global Optimization

(SMGO) algorithm for Lipschitz-continuous black-box functions is introduced. It employs the concepts of SM lower- and upper bounds in exploitation and exploration routines. Unlike the above-mentioned Lipschitz-based algorithms, we address the exploration problem systematically, by casting it as minimization of the uncertainty as inferred using SM techniques. Furthermore, our proposed algorithm is completely repeatable, and allows one to consider custom starting points and initial samples when available;

- Theoretical properties of SMGO are derived, including convergence, optimality gap, and computational complexity;
- Practical aspects are discussed, and ways to improve the computational efficiency are described. In particular, an iterative implementation is described, which exploits results from previous iterations to alleviate the computations at the current one;
- The SMGO algorithm is compared with other representative optimizers with different test functions, showing that the performance of SMGO is very competitive w.r.t. the state of the art, especially on functions with many global minima.

A preliminary version of this work appeared in [23]. With respect to that contribution, the algorithm presented here is more efficient, moreover the theoretical results and implementation analysis are new, and more benchmark results are presented. The MATLAB code for SMGO, as well as the other scripts used to generate the results, are available on GitHub under the following URL: <https://github.com/lorenzosabugjr/smgo>.

This paper is organized in seven sections. Section 2 gives the general problem statement and assumptions. Section 3 describes the proposed SMGO algorithm. The convergence properties and optimality gap calculations are discussed in Section 4, with implementation notes in Section 5. Section 6 compares the performance of the proposed algorithm with other global optimization techniques on representative test functions, and Section 7 concludes this paper and provides insights on future directions.

2. Problem Statement

Consider a cost function $z = f_o(\mathbf{x})$, $f_o : \mathcal{X} \rightarrow \mathbb{R}$, where $x \in \mathcal{X}$ is the vector of decision variables, $\mathcal{X} \subset \mathbb{R}^D$ is a compact and convex polytope (“search set”). No analytical form of f_o is assumed available. The *a priori* knowledge about f_o is given by the following assumption:

Assumption 1. f_o is Lipschitz continuous with unknown Lipschitz constant γ_o , i.e.,

$$f_o \in \mathcal{F}(\gamma_o)$$

where

$$\mathcal{F}(\gamma_o) \doteq \left\{ f \in \mathbf{C}^0(\mathcal{X}) : |f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq \gamma_o \|\mathbf{x}_1 - \mathbf{x}_2\|, \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X} \right\}$$

Being a Lipschitz continuous function on a compact set, f_o presents a global minimum z^* :

$$z^* = \min_{\mathbf{x} \in \mathcal{X}} f_o(\mathbf{x}). \quad (1)$$

Let us denote with

$$\mathcal{X}^* = \{ \mathbf{x} \in \mathcal{X} \mid f_o(\mathbf{x}) = z^* \} \quad (2)$$

the corresponding set of global minimizers.

We further assume that it is possible to acquire information about f_o by sampling (long function evaluation):

Assumption 2. Given a vector of decision variables $\mathbf{x} \in \mathcal{X}$, it is possible to sample the cost function without noise:

$$z = f_o(\mathbf{x}).$$

The set of collected data is denoted as:

$$\mathbf{X}^{(n)} = \left\{ (\mathbf{x}^{(1)}, z^{(1)}); (\mathbf{x}^{(2)}, z^{(2)}); \dots; (\mathbf{x}^{(n)}, z^{(n)}) \right\} \quad (3)$$

where $n \in \mathbb{N}$ is the number of data points and $z^{(i)} = f_o(\mathbf{x}^{(i)})$. For simplicity, with a slight abuse of notation, we will also write $\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}$ when the pair $(\mathbf{x}^{(i)}, z^{(i)})$ belongs to the data set (3).

We denote with $(\mathbf{x}^{*(n)}, z^{*(n)})$ the best pair in $\mathbf{X}^{(n)}$, where

$$(\mathbf{x}^{*(n)}, z^{*(n)}) = \underset{(\mathbf{x}^{(k)}, z^{(k)}) \in \mathbf{X}^{(n)}}{\arg \min} z^{(k)} \quad (4)$$

If the result of (4) is not unique, a lexicographic criterion is used to sort the minimizers, and the first one is picked. We further denote with $\delta^{(n)}$ the optimality gap:

$$\delta^{(n)} = z^{*(n)} - z^*. \quad (5)$$

We assume that a starting number $n_0 \geq 1$ of data points is also available, for example collected by the user in a first testing campaign, forming the set $\mathbf{X}^{(n_0)} \subset \mathcal{X}$.

We can now state the problem addressed in this paper.

Problem 1. *Design an algorithm that, under Assumptions 1-2, generates a sequence of points $\{\mathbf{x}^{(n_0+1)}, \mathbf{x}^{(n_0+2)}, \dots\}$, $\mathbf{x}^{(i)} \in \mathcal{X}$, such that:*

$$\forall \epsilon > 0, \exists n_\epsilon < \infty : z^{*(n_\epsilon)} \leq z^* + \epsilon. \quad (6)$$

Moreover, for a finite sequence length N , provide a method to compute a bound on the optimality gap $\delta^{(N)}$.

3. Set Membership Global Optimization (SMGO): Algorithm

SMGO addresses Problem 1 with a sequential procedure, as common in global optimization [1, 13, 18]. At each iteration $n \geq n_0$, the next test point $\mathbf{x}^{(n+1)}$ is chosen with a strategy that exploits the prior knowledge on function f_o , given by Assumption 1, and the current data set $\mathbf{X}^{(n)}$. In particular, a valid (i.e., consistent with data) estimate of the Lipschitz constant γ_o is derived at n_0 and updated at each $n > n_0$, in order to estimate upper and lower bounds on f_o . SMGO then leverages the latter to compute $\mathbf{x}^{(n+1)}$, choosing between an exploitation mode and an exploration one.

3.1. Preliminaries: Lipschitz constant estimation, cost function bounds, and search mesh

At the start of each iteration ($n > n_0$), one long function evaluation $z^{(n)} = f_o(\mathbf{x}^{(n)})$ is carried out, where $\mathbf{x}^{(n)}$ is the test point selected at the previous iteration. This resulting new data point $(\mathbf{x}^{(n)}, z^{(n)})$ is added to the set of collected samples:

$$\mathbf{X}^{(n)} = \mathbf{X}^{(n-1)} \cup (\mathbf{x}^{(n)}, z^{(n)}).$$

At the first iteration, $n = n_0$, an initial estimate of the Lipschitz constant, $\gamma^{(n_0)}$, is computed as:

$$\gamma^{(n_0)} = \max_{i,j=1,\dots,n_0, i \neq j} \frac{|z^{(i)} - z^{(j)}|}{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|} \quad (7)$$

If $n_0 = 1$, the Lipschitz constant estimate can be initialized as an arbitrarily small number. When $n > n_0$, the minimum feasible Lipschitz constant estimate $\gamma^{(n)}$ is updated as follows:

$$\gamma^{(n)} = \max \left(\gamma^{(n-1)}, \max_{k=1,\dots,n-1} \frac{|z^{(n)} - z^{(k)}|}{\|\mathbf{x}^{(n)} - \mathbf{x}^{(k)}\|} \right). \quad (8)$$

Such an estimate is guaranteed to be compatible with the available data and prior information on f_o , given Assumption 2 [20]. Moreover, $\gamma^{(n)}$ is monotonically increasing with n and convergence to the true Lipschitz constant γ_o is guaranteed if the test points densely cover the search set [22]. This is the case for the SMGO algorithm as shown in Section 4.

Given the data set $\mathbf{X}^{(n)}$ and Lipschitz constant estimate $\gamma^{(n)}$ we can compute the following lower and upper bounds on f_o [20]:

$$\underline{z}^{(n)}(\mathbf{x}) \triangleq \max_{k=1,\dots,n} \left(z^{(k)} - \mu \gamma^{(n)} \|\mathbf{x} - \mathbf{x}^{(k)}\| \right), \quad (9)$$

$$\bar{z}^{(n)}(\mathbf{x}) \triangleq \min_{k=1,\dots,n} \left(z^{(k)} + \mu \gamma^{(n)} \|\mathbf{x} - \mathbf{x}^{(k)}\| \right). \quad (10)$$

When $\mu = 1$, the functions $\underline{z}^{(n)}$, $\bar{z}^{(n)}$ represent, respectively, the minimum and maximum values that the objective function can take for a given vector of decision variables $\mathbf{x} \in \mathcal{X}$, given the information collected up to iteration n .

In SMGO, an (user-defined) overestimation parameter $\mu > 1$ is used in (9)-(10) to compensate for the fact that the $\gamma^{(n)}$ estimated from the currently available data points is an underestimator of the true Lipschitz constant γ_o .

The set

$$\mathbb{F}^{(n)} = \left\{ f \in \mathcal{F}(\gamma^{(n)}) : \underline{z}^{(n)}(\mathbf{x}) \leq f(\mathbf{x}) \leq \bar{z}^{(n)}(\mathbf{x}) \right\}$$

is referred to as the *feasible functions set* [20] at iteration n . The functions $\underline{z}^{(n)}(\mathbf{x})$ and $\bar{z}^{(n)}(\mathbf{x})$ are intersections of hypercones (referred simply as cones

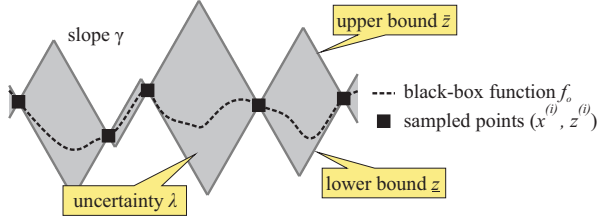


Figure 1: Scalar example of optimal lower and upper bounds $\underline{z}^{(n)}(\mathbf{x})$, $\bar{z}^{(n)}(\mathbf{x})$ from finite samples of f_o , and of the corresponding uncertainty interval $\lambda^{(n)}(\mathbf{x})$.

in the remainder), each one with vertex at a tested point. A representation of these bounds is depicted in Fig. 1 for a function f_o with scalar input argument. It should be emphasized that $\underline{z}^{(n)}$ and $\bar{z}^{(n)}$ are calculated from the data in $\mathbf{X}^{(n)}$, without any need for additional evaluations of f_o .

The uncertainty about the value of f_o at the point $\mathbf{x} \in \mathcal{X}$ is computed as:

$$\lambda^{(n)}(\mathbf{x}) = \bar{z}^{(n)}(\mathbf{x}) - \underline{z}^{(n)}(\mathbf{x}). \quad (11)$$

The SMGO algorithm proceeds in two steps. First, an exploitation routine searches for a local minimizer of $\underline{z}^{(n)}$. Then, an exploration step seeks a maximizer of $\lambda^{(n)}$. Exact solutions for such subproblems can be computed by utilizing Hyperbolic Voronoi Diagrams (HVDs) [20, 24], however at rather high computational cost. To increase computational efficiency, we thus introduce a mesh given by segments $l_{i,j}$ connecting any two points $\mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in \mathbf{X}^{(n)}$:

$$l_{i,j} = \left\{ \mathbf{x} = a\mathbf{x}^{(i)} + (1-a)\mathbf{x}^{(j)} : a \in [0,1], \right. \\ \left. \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in \mathbf{X}^{(n)}, i \neq j \right\} \quad (12)$$

Such a mesh allows us to derive analytic solutions to the computation of the next test point $\mathbf{x}^{(n+1)}$, according to the two strategies presented next.

3.2. Exploitation (Mode θ)

Consider the best data point $(\mathbf{x}^{*(n)}, z^{*(n)})$ (4) at iteration n . Mode θ searches a point $\mathbf{x}_{\theta^*}^{(n)}$ with the highest predicted improvement w.r.t. $\mathbf{x}^{*(n)}$, according to the lower bound $\underline{z}^{(n)}$ (9). Specifically, the exploitation routine considers the set of segments $l_{n^*,i}$ between $\mathbf{x}^{*(n)}$ and any other sampled point $\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}$ to search for the candidate point

$\mathbf{x}_{\theta^*}^{(n)}$. Here, index n^* is the data-point index of the pair $(\mathbf{x}^{*(n)}, z^{*(n)})$ in the data set $\mathbf{X}^{(n)}$.

We first evaluate the lower bounds on $f_o(\mathbf{x}) : \mathbf{x} \in l_{n^*,i}$ considering only the hypercones defined by the data points $(\mathbf{x}^{*(n)}, z^{*(n)})$ and $(\mathbf{x}^{(i)}, z^{(i)})$, i.e.:

$$z^{*(n)} - \mu\gamma^{(n)}\|\mathbf{x} - \mathbf{x}^{*(n)}\| \quad (13)$$

and

$$z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x} - \mathbf{x}^{(i)}\|. \quad (14)$$

The intersection of these two cones on the segment $l_{n^*,i}$ is the smallest feasible value of $f_o(\mathbf{x}) : \mathbf{x} \in l_{n^*,i}$ considering the information given by the two points $(\mathbf{x}^{*(n)}, z^{*(n)})$, $(\mathbf{x}^{(i)}, z^{(i)})$ and the Lipschitz constant estimate $\gamma^{(n)}$. Such an intersection, denoted with $\mathbf{x}_{\theta}^{(i)}$, can be derived analytically on the basis of geometrical considerations (see Fig. 2). In fact, denoting

$$s_i = \frac{z^{(i)} - z^{*(n)}}{\|\mathbf{x}^{(i)} - \mathbf{x}^{*(n)}\|} \geq 0, \quad (15)$$

it follows that

$$\mathbf{x}_{\theta}^{(i)} = \mathbf{x}^{*(n)} + \frac{1 - \frac{s_i}{\mu\gamma^{(n)}}}{2}(\mathbf{x}^{(i)} - \mathbf{x}^{*(n)}). \quad (16)$$

By evaluating (16) for all possible segments, we obtain the set $\mathbf{X}_{\theta}^{(n)} = \{\mathbf{x}_{\theta}^{(i)}, i = 1, \dots, n, i \neq n^*\}$. Then $\mathbf{x}_{\theta^*}^{(n)}$ is chosen as:

$$\mathbf{x}_{\theta^*}^{(n)} = \arg \min_{\mathbf{x} \in \mathbf{X}_{\theta}^{(n)}} \underline{z}^{(n)}(\mathbf{x}) \quad (17a)$$

s.t.

$$\underline{z}^{(n)}(\mathbf{x}) = z^{*(n)} - \mu\gamma^{(n)}\|\mathbf{x} - \mathbf{x}^{*(n)}\| \quad (17b)$$

where the constraint (17b) ensures that the point selected is such that the lower bound $\underline{z}^{(n)}$ considering the information provided by all the sampled points coincides with the lower bound provided by the current best point $(\mathbf{x}^{*(n)}, z^{*(n)})$ alone. This usually happens in the vicinity of $\mathbf{x}^{*(n)}$, in particular at the border of the hyperbolic Voronoi cell pertaining to $\mathbf{x}^{*(n)}$, see [20]. The feasibility of problem (17) is always guaranteed, as shown in Lemma 1 in Section 4.

To decide whether $\mathbf{x}_{\theta^*}^{(n)}$ shall be eventually taken as the next test point $\mathbf{x}^{(n+1)}$ for the long function, or the exploration routine shall be called instead, the following condition is evaluated:

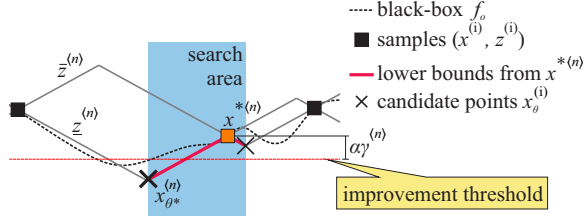


Figure 2: Qualitative example of the exploitation routine for a scalar problem.

$$\underline{z}^{(n)}(\mathbf{x}_{\theta^*}^{(n)}) \leq z^{*(n)} - \alpha\gamma^{(n)}, \quad (18)$$

where $\alpha\gamma^{(n)}$ is referred to as the *expected improvement threshold* and $\alpha \in [0, 1]$ is the SMGO *exploitation parameter*, adjusted by the user. If (18) is fulfilled (as in the example of Fig. 2), it means that the potential improvement provided by $\mathbf{x}_{\theta^*}^{(n)}$ is larger than the threshold. In this case, we set

$$\mathbf{x}^{(n+1)} = \mathbf{x}_{\theta^*}^{(n)}.$$

Otherwise, SMGO switches to exploration mode, described in the following subsection.

3.3. Exploration (Mode ψ)

In this mode, the candidate point $\mathbf{x}_{\psi^*}^{(n)} \in \mathcal{X}$ is chosen as the one with largest $\lambda^{(n)}$ value among selected candidates. Sampling $\mathbf{x}_{\psi^*}^{(n)}$ will yield $\lambda^{(n+1)}(\mathbf{x}_{\psi^*}^{(n)}) = 0$ (due to noiseless function evaluation, see Assumption 2), and the uncertainty in its vicinity will correspondingly decrease.

Similar to Mode θ , the search for $\mathbf{x}_{\psi^*}^{(n)}$ considers the segments among the samples in $\mathbf{X}^{(n)}$, in order to limit the computational complexity. In this case, we consider as candidates all the midpoints between any pair $(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \in \mathbf{X}^{(n)}$. The reason is that midpoints generally feature larger uncertainty than other points on the segments, which are closer to one of the two extremes. Hence, we form the set of midpoints:

$$\mathbf{X}_{\psi}^{(n)} = \left\{ \mathbf{x}_{\psi}^{(i,j)} = \frac{\mathbf{x}^{(i)} + \mathbf{x}^{(j)}}{2} \mid \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in \mathbf{X}^{(n)} \right\}, \quad (19)$$

and, according to the described rationale, the exploration routine picks the next test point as:

$$\mathbf{x}^{(n+1)} = \mathbf{x}_{\psi^*}^{(n)} = \arg \max_{\mathbf{x} \in \mathbf{X}_{\psi}^{(n)}} \lambda^{(n)}(\mathbf{x}). \quad (20)$$

Also in this case, if (20) does not have a unique solution, a lexicographic criterion is added to select the maximizer.

3.4. Algorithm summary

A summarized flow of the SMGO method is given as pseudo-code in Algorithm 1.

Algorithm 1: SMGO Algorithm

Input: Long function f_o , initial data points $\mathbf{X}^{(n_0)}$, maximum number of long function evaluations N , exploitation parameter α , overestimation parameter μ

- 1 **while** $n_0 < n \leq N$ **do**
 - // Long function evaluation, data update
 - 2 Evaluate the long function f_o at $\mathbf{x}^{(n)}$ and measure output $z^{(n)}$. Update the data set $\mathbf{X}^{(n)} \leftarrow \mathbf{X}^{(n-1)} \cup (\mathbf{x}^{(n)}, z^{(n)})$
 - 3 Update current best sample $(\mathbf{x}^{*(n)}, z^{*(n)})$ and Lipschitz constant $\gamma^{(n)}$ from $\mathbf{X}^{(n)}$
 - // Exploitation (Mode θ)
 - 4 Evaluate lower bounds on segments around the current best sample $\mathbf{x}^{*(n)}$ to choose the exploitation point $\mathbf{x}_{\theta^*}^{(n)}$
 - 5 **if** *expected improvement condition* (18) **is met then**
 - 6 Assign test point for next iteration $\mathbf{x}^{(n+1)} \leftarrow \mathbf{x}_{\theta^*}^{(n)}$
 - 7 **else**
 - // Exploration (Mode ψ)
 - 8 Find $\mathbf{x}_{\psi^*}^{(n)}$ as the midpoint, over all possible segments, with largest uncertainty $\lambda^{(n)}$
 - 9 Assign test point for next iteration $\mathbf{x}^{(n+1)} \leftarrow \mathbf{x}_{\psi^*}^{(n)}$
 - 10 Go to next iteration $n \leftarrow n + 1$
- 11 Final optimal point and value: get the best sample $(\mathbf{x}^{*(N)}, z^{*(N)})$ from the set $\mathbf{X}^{(N)}$

3.5. On the choice of α and μ

The SMGO parameter α affects only Mode θ : a smaller α value leads to higher tendency for exploitation. Regarding μ , using a value close to one is recommended, because a higher Lipschitz constant overestimator leads to $\mathbf{x}_{\theta}^{(i)} \rightarrow \frac{\mathbf{x}^{*(n)} + \mathbf{x}^{(i)}}{2}$, which would not fully utilize the Set Membership

approach in exploitation. At the same time, excessively large values of μ increase the conservativeness of $\lambda^{(n)}$. Moreover, in the computation of $\lambda^{(n)}(\mathbf{x}_\psi^{(i,j)})$, $\mathbf{x}_\psi^{(i,j)} \in \mathbf{X}_\psi^{(n)}$ (19)-(20), with larger μ values the terms $\|\mathbf{x}_\psi^{(i,j)} - \mathbf{x}^{(i)}\|$, $i, j = 1, \dots, n$ gain a higher relative importance with respect to the corresponding values of $z^{(i)}$ (see (9)-(10)), so that the exploration mode will tend to select the mid-points that are farthest away from the available data-points.

4. Theoretical Properties

We now analyze the convergence and optimality properties of the SMGO algorithm. First, we prove that problem (17) is always feasible, thus showing that the exploitation routine always returns a valid candidate next point $\mathbf{x}_{\theta^*}^{(n)}$.

Lemma 1. *At any iteration n , problem (17) admits at least one feasible point.*

Proof. Take the pair $(\mathbf{x}^{(i)}, z^{(i)})$ such that

$$z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}^{(i)}\| = \max_{(\mathbf{x}, z) \in \mathbf{X}^{(n)} \setminus \{\mathbf{x}^{*(n)}, z^{*(n)}\}} z - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}\|, \quad (21)$$

and consider the corresponding candidate exploitation point $\mathbf{x}_\theta^{(i)}$ given by (16). By construction we have that

$$z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(i)}\| = z^{*(n)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{*(n)}\|,$$

and we want now to prove that this point also coincides with $\underline{z}^{(n)}(\mathbf{x}_\theta^{(i)})$, i.e. that:

$$z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(i)}\| = \max_{(\mathbf{x}, z) \in \mathbf{X}^{(n)}} z - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}\|. \quad (22)$$

Assume now, for the purpose of contradiction, that:

$$\begin{aligned} & \exists (\mathbf{x}^{(j)}, z^{(j)}) \neq (\mathbf{x}^{(i)}, z^{(i)}) : \\ & z^{(j)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(j)}\| > z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(i)}\|, \end{aligned} \quad (23)$$

thus invalidating (22). Then, denoting with $a = \frac{1}{2}(1 - \frac{s_i}{\mu\gamma^{(n)}})$ (see (16)), we would have:

$$\begin{aligned} & z^{(j)} - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}^{(j)}\| \geq \\ & z^{(j)} - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}_\theta^{(i)}\| - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(j)}\| > \\ & z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(i)}\| - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}_\theta^{(i)}\| = \\ & z^{(i)} - \mu\gamma^{(n)}\left(\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(i)}\| + \|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{*(n)}\|\right) = \\ & z^{(i)} - \mu\gamma^{(n)}\left((1-a)\|\mathbf{x}^{*(n)} - \mathbf{x}^{(i)}\| + a\|\mathbf{x}^{*(n)} - \mathbf{x}^{(i)}\|\right) = \\ & z^{(i)} - \mu\gamma^{(n)}\|\mathbf{x}^{*(n)} - \mathbf{x}^{(i)}\|, \end{aligned}$$

which contradicts equation (21). Then, (23) is false, while (22) holds true, implying that point $\mathbf{x}_\theta^{(i)}$ satisfies constraint (17b), hence proving the result. \square

4.1. Convergence of SMGO algorithm

We next show that after finite iterations, an ϵ -suboptimal point is obtained, as required by Problem 1. For simplicity, we consider the following technical assumption, which in practice can be replaced by a vertex estimation mechanism presented in Section 5.1. Let us denote with $V \in \mathbb{N}$ the number of vertices of the polytope \mathcal{X} , and with $\mathbf{c}^{(v)}$ the v -th vertex.

Assumption 3. *The starting data set $\mathbf{X}^{(n_0)}$ contains all the vertices of \mathcal{X} : $\mathbf{c}^{(v)} \in \mathbf{X}^{(n_0)}$, $\forall v \in [1, \dots, V]$.*

We start by proving four lemmas that are instrumental to show SMGO convergence.

In the following, for a point \mathbf{x} and positive scalar r , denote with $\mathcal{B}(\mathbf{x}, r)$ the closed hyper-ball:

$$\mathcal{B}(\mathbf{x}, r) = \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| \leq r\}$$

Lemma 2. *Mode θ will fail after a finite number of iterations.*

Proof. At iteration n , denote with $\mathcal{G}^{(n)}$ the subset of the search space where a candidate test point is accepted to perform an exploitation, i.e.,

$$\mathcal{G}^{(n)} \triangleq \left\{ \mathbf{x} \in \mathcal{X} : \underline{z}^{(n)}(\mathbf{x}) < z^{*(n)} - \alpha\gamma^{(n)} \right\}.$$

When a candidate point inside this set is chosen as $\mathbf{x}^{(n+1)}$ in Mode θ , two cases arise:

1. The new sample does not change the current best point $\mathbf{x}^{*(n)}$. In this case, $z^{(n+1)} \geq z^{*(n)}$.

Denote as $\mathcal{B}(\mathbf{x}^{(n+1)}, r_\theta)$, the hyper-ball centered at $\mathbf{x}^{(n+1)}$ with finite radius

$$r_\theta \triangleq \frac{z^{(n+1)} - (z^{*(n)} - \alpha\gamma^{(n)})}{\mu\gamma^{(n+1)}}.$$

It follows that $\forall \mathbf{x} \in \mathcal{B}(\mathbf{x}^{(n+1)}, r_\theta)$

$$\underline{z}^{(n+1)}(\mathbf{x}) \geq z^{*(n)} - \alpha\gamma^{(n)}.$$

Therefore all the points inside the hyper-ball are not eligible for exploitation. Hence, we have:

$$\mathcal{G}^{(n+1)} = \mathcal{G}^{(n)} \setminus \mathcal{B}(\mathbf{x}^{(n+1)}, r_\theta)$$

i.e., \mathcal{G} diminishes by a finite amount, which is valid even when $\gamma^{(n)}$ updates, because $\gamma^{(n)}$ is bounded by γ_o .

2. The new sample replaces the current best point $\mathbf{x}^{*(n)}$. In this case $z^{(n+1)} < z^{*(n)}$ and therefore $\mathcal{G}^{(n+1)} \subset \mathcal{G}^{(n)}$ due to the new threshold. Moreover, the hyper-ball $\mathcal{B}(\mathbf{x}^{(n+1)}, r_\theta)$, with $r_\theta = \frac{\alpha}{\mu}$ around the new sample is also removed from \mathcal{G} .

Hence, $\mathcal{G}^{(n)} = \emptyset$ after finite iterations, and Mode θ will fail, proving the lemma. \square

Lemma 3. *Let Assumption 3 hold and consider any unsampled point $\hat{\mathbf{x}}$ in the interior of \mathcal{X} and any sample $\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}$. Consider the open half-space:*

$$\mathcal{Q} = \left\{ \mathbf{x} \in \mathcal{X} : (\mathbf{x} - \mathbf{x}^{(i)})^\top (\hat{\mathbf{x}} - \mathbf{x}^{(i)}) > 0 \right\}.$$

Then, there exists at least one sample $\mathbf{x}^{(j)} \in \mathbf{X}^{(n)} \cap \mathcal{Q}$, $i \neq j$.

Proof. Let us assume, for the purpose of contradiction, that the lemma does not hold. That is, there is no sample in the open half-space \mathcal{Q} . Due to Assumption 3, this would also apply to all vertices $\mathbf{c}^{(v)}$, $v = 1, \dots, V$ of \mathcal{X} . Therefore, we would have:

$$(\hat{\mathbf{x}} - \mathbf{x}^{(i)})^\top (\mathbf{c}^{(v)} - \mathbf{x}^{(i)}) \leq 0, \forall v = 1, \dots, V. \quad (24)$$

Two cases could then arise:

1. Inequality in (24) applies strictly for at least one vertex $\mathbf{c}^{(v)}$. In this case, any such vertex $\mathbf{c}^{(v)}$ and $\hat{\mathbf{x}}$ would reside on opposite sides of the hyper-plane

$$\mathcal{P} \triangleq \left\{ \mathbf{x} : (\hat{\mathbf{x}} - \mathbf{x}^{(i)})^\top (\mathbf{x} - \mathbf{x}^{(i)}) = 0 \right\}.$$

However, consider any $\mathbf{w} \in \mathcal{X}$. Since the latter is a convex polytope, we can write

$$\mathbf{w} = \sum_{v=1}^V b_v \mathbf{c}^{(v)} \quad (25)$$

with $0 \leq b_v \leq 1$, $v = 1, \dots, V$, and $\sum_{v=1}^V b_v = 1$.

Hence, it would apply that

$$\begin{aligned} \forall \mathbf{w} \in \mathcal{X}, (\hat{\mathbf{x}} - \mathbf{x}^{(i)})^\top (\mathbf{w} - \mathbf{x}^{(i)}) &= \\ (\hat{\mathbf{x}} - \mathbf{x}^{(i)})^\top \left(\sum_{v=1}^V b_v \mathbf{c}^{(v)} - \sum_{v=1}^V b_v \mathbf{x}^{(i)} \right) &= \\ \sum_{v=1}^V b_v (\hat{\mathbf{x}} - \mathbf{x}^{(i)})^\top (\mathbf{c}^{(v)} - \mathbf{x}^{(i)}) &< 0, \end{aligned}$$

which would mean that all $\mathbf{w} \in \mathcal{X}$ reside on the opposite side of \mathcal{P} with respect to $\hat{\mathbf{x}}$. This can only happen if $\hat{\mathbf{x}} \notin \mathcal{X}$, falling in a contradiction with the assumptions.

2. Equality in (24) applies for all $\mathbf{c}^{(v)}$, $v = 1, \dots, V$. This would mean that all vertices (and consequently, all $\mathbf{w} \in \mathcal{X}$) belong to the hyper-plane \mathcal{P} containing $\mathbf{x}^{(i)}$. In turn, this would imply that either $\hat{\mathbf{x}} \notin \mathcal{X}$ (if $\|\hat{\mathbf{x}} - \mathbf{x}^{(i)}\| > 0$), or $\hat{\mathbf{x}} = \mathbf{x}^{(i)}$, both of which fall in contradiction with the assumptions.

Since both cases 1. and 2. lead to contradiction, the claim of the lemma must be true, thus completing the proof. \square

Lemma 4. *Consider any unsampled point $\hat{\mathbf{x}}$ and any radius r such that $\mathbf{X}^{(n)} \cap \mathcal{B}(\hat{\mathbf{x}}, r) = \emptyset$. Then, given Assumption 2 and a Lipschitz constant estimate $\gamma^{(n)}$ obtained with (7)-(8), it applies that*

$$\lambda^{(n)}(\hat{\mathbf{x}}) \geq 2(\mu - 1)\gamma^{(n)}r. \quad (26)$$

Proof. Consider the point $\hat{\mathbf{x}}$ and denote with $(\mathbf{x}^{(a)}, z^{(a)})$, $(\mathbf{x}^{(b)}, z^{(b)}) \in \mathbf{X}^{(n)}$ the data pairs that determine the upper and lower bounds $\bar{z}^{(n)}(\hat{\mathbf{x}})$, $\underline{z}^{(n)}(\hat{\mathbf{x}})$, respectively, i.e:

$$\begin{aligned} \bar{z}^{(n)}(\hat{\mathbf{x}}) &= \min_{k=1 \dots n} \left(z^{(k)} + \mu\gamma^{(n)}\|\hat{\mathbf{x}} - \mathbf{x}^{(k)}\| \right) = \\ & z^{(a)} + \mu\gamma^{(n)}\|\hat{\mathbf{x}} - \mathbf{x}^{(a)}\| \end{aligned}$$

$$\underline{z}^{(n)}(\hat{\mathbf{x}}) = \max_{k=1,\dots,n} \left(z^{(k)} - \mu\gamma^{(n)} \|\hat{\mathbf{x}} - \mathbf{x}^{(k)}\| \right) = z^{(b)} - \mu\gamma^{(n)} \|\hat{\mathbf{x}} - \mathbf{x}^{(b)}\|.$$

Then, considering that $\mathbf{X}^{(n)} \cap \mathcal{B}(\hat{\mathbf{x}}, r) = \emptyset$, that $(\mu - 1) > 0$, and taking into account (7)-(8) we have:

$$\begin{aligned} \lambda^{(n)}(\hat{\mathbf{x}}) &= \bar{z}^{(n)}(\hat{\mathbf{x}}) - \underline{z}^{(n)}(\hat{\mathbf{x}}) \\ &= z^{(a)} - z^{(b)} + \mu\gamma^{(n)} \left(\|\hat{\mathbf{x}} - \mathbf{x}^{(a)}\| + \|\hat{\mathbf{x}} - \mathbf{x}^{(b)}\| \right) \\ &\geq -\gamma^{(n)} \|\mathbf{x}^{(a)} - \mathbf{x}^{(b)}\| + \mu\gamma^{(n)} \left(\|\hat{\mathbf{x}} - \mathbf{x}^{(a)}\| + \|\hat{\mathbf{x}} - \mathbf{x}^{(b)}\| \right) \\ &\geq -\gamma^{(n)} \left(\|\hat{\mathbf{x}} - \mathbf{x}^{(a)}\| + \|\hat{\mathbf{x}} - \mathbf{x}^{(b)}\| \right) + \mu\gamma^{(n)} \left(\|\hat{\mathbf{x}} - \mathbf{x}^{(a)}\| + \|\hat{\mathbf{x}} - \mathbf{x}^{(b)}\| \right) \\ &\geq 2(\mu - 1)\gamma^{(n)}r \end{aligned}$$

which proves the result. \square

Note that in Lemma 4 we are not assuming that the estimate $\gamma^{(n)}$ is larger than γ_o : the result follows by how the estimate is computed, combined with the use of $\mu > 1$.

Lemma 5. *Let Assumption 3 hold, and assume that Mode ψ is undertaken infinitely often as $n \rightarrow +\infty$. Then, for any $\hat{\mathbf{x}} \in \mathcal{X}$ and any $\sigma > 0$, $\exists n_\sigma < \infty$ such that*

$$\min_{\mathbf{x}^{(i)} \in \mathbf{X}^{(n_\sigma)}} \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}\| < \sigma.$$

Proof. Consider any point $\hat{\mathbf{x}} \in \mathcal{X}$. If $\hat{\mathbf{x}} \in \mathbf{X}^{(n)}$ for some $n < \infty$, we have $\min_{\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}} \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}\| = 0$ and the claim is trivially proven. Consider then the case $\hat{\mathbf{x}} \notin \mathbf{X}^{(n)}, \forall n \in [1, +\infty)$. Pick the nearest sample

$$\bar{\mathbf{x}}^{(n)} = \min_{\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}} \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}\| \quad (27)$$

and consider the open half-space:

$$\mathcal{Q} = \left\{ \mathbf{x} \in \mathcal{X} : (\mathbf{x} - \bar{\mathbf{x}}^{(n)})^\top (\hat{\mathbf{x}} - \bar{\mathbf{x}}^{(n)}) > 0 \right\}.$$

Now, we select sample $\tilde{\mathbf{x}}^{(n)}$

$$\tilde{\mathbf{x}}^{(n)} = \arg \min_{\mathbf{x}^{(j)} \in \mathbf{X}^{(n)} \cap \mathcal{Q}} \|\mathbf{x}^{(j)} - \bar{\mathbf{x}}^{(n)}\|. \quad (28)$$

Such a sample is always guaranteed to exist, in force of Lemma 3. We note that due to our selection criterion for $\tilde{\mathbf{x}}^{(n)}$, the set

$$\mathcal{H} = \left\{ \mathbf{x} \in \mathcal{Q} : \|\mathbf{x} - \bar{\mathbf{x}}^{(n)}\| < \|\tilde{\mathbf{x}}^{(n)} - \bar{\mathbf{x}}^{(n)}\| \right\}$$

does not have any sample in its interior, otherwise any such sample should have been selected as $\tilde{\mathbf{x}}^{(n)}$ in (28). The geometry of \mathcal{H} is an open half-hyperball, with center at $\bar{\mathbf{x}}^{(n)}$ and radius $\|\tilde{\mathbf{x}}^{(n)} - \bar{\mathbf{x}}^{(n)}\|$.

We now show that a midpoint inside \mathcal{H} is going to be sampled after finite iterations. There exists at least the midpoint $\hat{\mathbf{m}} \triangleq \frac{\bar{\mathbf{x}}^{(n)} + \tilde{\mathbf{x}}^{(n)}}{2}$ in the interior of \mathcal{H} . Consider the ball $\mathcal{B}(\hat{\mathbf{m}}, h)$ centered at $\hat{\mathbf{m}}$ with radius h :

$$\begin{aligned} h &= \max_{a \in \mathbb{R}^+} a \\ &\text{s.t.} \\ \mathcal{B}(\hat{\mathbf{m}}, a) &\subset \mathcal{H} \end{aligned} \quad (29)$$

Then, recalling that there are no samples inside \mathcal{H} , by Lemma 4 we have:

$$\lambda^{(n)}(\hat{\mathbf{m}}) \geq 2(\mu - 1)\gamma^{(n)}h. \quad (30)$$

Furthermore, we consider the set of points not belonging to \mathcal{H} , whose uncertainty is larger than that of $\hat{\mathbf{m}}$

$$\mathcal{R}^{(n)}(\hat{\mathbf{m}}) = \left\{ \mathbf{x} \notin \mathcal{H} : \lambda^{(n)}(\mathbf{x}) > \lambda^{(n)}(\hat{\mathbf{m}}) \right\}.$$

We then sort all candidate midpoints in the set $\mathbf{X}_\psi^{(n)}$, by decreasing $\lambda^{(n)}$ values:

$$\mathbf{M}^{(n)} \triangleq \left\{ \mathbf{m}_1^{(n)}, \dots, \mathbf{m}_M^{(n)} \right\} : \lambda^{(n)}(\mathbf{m}_k^{(n)}) \geq \lambda^{(n)}(\mathbf{m}_{k+1}^{(n)})$$

where $\mathbf{m}_k^{(n)} = \mathbf{x}_{\psi}^{(i,j)}$ for some $\mathbf{x}^{(i)}, \mathbf{x}^{(j)} \in \mathbf{X}^{(n)}$ (see (19)) and $M = \frac{n(n-1)}{2}$ is the total number of segments at iteration n . In Mode ψ , the first-ranking candidate $\mathbf{m}_1^{(n)}$ is taken for sampling. Its uncertainty is necessarily larger than $\lambda^{(n)}(\hat{\mathbf{m}})$ (otherwise $\hat{\mathbf{m}}$ would have been ranked higher), i.e. $\mathbf{m}_1^{(n)} \in \mathcal{R}^{(n)}(\hat{\mathbf{m}})$. Assume, for the purpose of contradiction, that points in \mathcal{H} are never sampled.

At iteration $n+1$, the midpoint $\mathbf{m}_1^{(n)} \notin \mathcal{H}$ is thus sampled and its uncertainty $\lambda^{(n+1)}(\mathbf{m}_1^{(n)})$ becomes zero. Since this point does not belong to \mathcal{H} , the value of h in (29) is still valid, while the uncertainty bound (30) may either be the same, or increase in case $\gamma^{(n+1)} > \gamma^{(n)}$ (see (8)). Moreover, for any

point $\mathbf{x} \in \mathcal{X}$ we have

$$\lambda^{(n+1)}(\mathbf{x}) \leq 2\mu\gamma^{(n+1)}\|\mathbf{x} - \mathbf{m}_1^{(n)}\| \quad (31)$$

as shown by direct application of (9)-(11). Consider now the ball $\mathcal{B}(\mathbf{m}_1^{(n)}, r_\lambda)$, where

$$r_\lambda = \left(1 - \frac{1}{\mu}\right)h. \quad (32)$$

On the basis of (30), (31) and (32) we have that:

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{B}(\mathbf{m}_1^{(n)}, r_\lambda), \lambda^{(n+1)}(\mathbf{x}) &\leq 2\mu\gamma^{(n+1)}r_\lambda = \\ 2(\mu - 1)\gamma^{(n+1)}h &\leq \lambda^{(n+1)}(\hat{\mathbf{m}}). \end{aligned} \quad (33)$$

Hence,

$$\mathcal{R}^{(n+1)}(\hat{\mathbf{m}}) = \mathcal{R}^{(n)}(\hat{\mathbf{m}}) \setminus \left(\mathcal{R}^{(n)}(\hat{\mathbf{m}}) \cap \mathcal{B}(\mathbf{m}_1^{(n)}, r_\lambda) \right).$$

This means that the volume of set $\mathcal{R}^{(n)}(\hat{\mathbf{m}})$ diminishes by a finite quantity at each iteration, since the value of h does not depend on n ; and becomes null after finite iterations, since Mode ψ is assumed to be undertaken infinitely often. Denote with $\bar{n} - 1 < \infty$ the iteration at which this happens. Then, no midpoint outside the set \mathcal{H} would have uncertainty larger than $\lambda^{(\bar{n})}(\hat{\mathbf{m}})$, so that a midpoint

$$\left\{ \mathbf{x}_\psi^{(i,j)} \in \mathbf{X}_\psi^{(\bar{n})} \cap \mathcal{H} : \lambda^{(\bar{n})}(\mathbf{x}_\psi^{(i,j)}) \geq \lambda^{(\bar{n})}(\hat{\mathbf{m}}) \right\}$$

will be ranked first in $\mathcal{M}^{(\bar{n}-1)}$ and sampled at iteration \bar{n} , thus falling in contradiction.

We thus demonstrated that, in finite iterations, a point inside \mathcal{H} is sampled. Moreover, denoting such a point as $\mathbf{x}^{(\bar{n})}$, by construction we have

$$\|\mathbf{x}^{(\bar{n})} - \bar{\mathbf{x}}^{(n)}\| < \|\tilde{\mathbf{x}}^{(n)} - \bar{\mathbf{x}}^{(n)}\|. \quad (34)$$

Then, two cases may occur: either $\mathbf{x}^{(\bar{n})}$ is closer to $\hat{\mathbf{x}}$ than $\bar{\mathbf{x}}^{(n)}$, or not. In the former case, repeating the process from (27) we'll have that $\bar{\mathbf{x}}^{(\bar{n})} = \mathbf{x}^{(\bar{n})}$ and $\tilde{\mathbf{x}}^{(\bar{n})} = \bar{\mathbf{x}}^{(n)}$. In the latter case, $\bar{\mathbf{x}}^{(\bar{n})} = \bar{\mathbf{x}}^{(n)}$ and $\tilde{\mathbf{x}}^{(\bar{n})} = \mathbf{x}^{(\bar{n})}$. In both cases, the radius of the half-hyper-ball \mathcal{H} shrinks because of (34). By induction we thus have:

$$\lim_{n \rightarrow \infty} \|\tilde{\mathbf{x}}^{(n)} - \bar{\mathbf{x}}^{(n)}\| = 0. \quad (35)$$

Next, still considering $\bar{\mathbf{x}}^{(n)}$ as the nearest sample to $\hat{\mathbf{x}}$ (27), we are going to show that there exists a

radius $a > 0$ such that

$$\|\mathbf{x} - \hat{\mathbf{x}}\| < \|\bar{\mathbf{x}}^{(n)} - \hat{\mathbf{x}}\|, \forall \mathbf{x} \in \mathcal{Q} \cap \mathcal{B}(\bar{\mathbf{x}}^{(n)}, a) \quad (36)$$

To this end, consider any point $\mathbf{x} \in \mathcal{Q}$ and assume, for the sake of contradiction, that (36) never holds. This would mean that the nearest point to $\hat{\mathbf{x}}$ on the segment connecting $\bar{\mathbf{x}}^{(n)}$ to \mathbf{x} is always $\bar{\mathbf{x}}^{(n)}$, no matter how close the two points are. In turn, this would imply that:

$$\frac{d}{da} \left\| ((1-a)\bar{\mathbf{x}}^{(n)} + a\mathbf{x}) - \hat{\mathbf{x}} \right\|_{a=0} \geq 0, \quad (37)$$

i.e.,

$$\begin{aligned} \frac{1}{\|\bar{\mathbf{x}}^{(n)} - \hat{\mathbf{x}}\|} (\bar{\mathbf{x}}^{(n)} - \hat{\mathbf{x}})^\top (\mathbf{x} - \bar{\mathbf{x}}^{(n)}) &\geq 0 \\ (\hat{\mathbf{x}} - \bar{\mathbf{x}}^{(n)})^\top (\mathbf{x} - \bar{\mathbf{x}}^{(n)}) &\leq 0 \end{aligned}$$

which would imply that $\mathbf{x} \notin \mathcal{Q}$, thus falling in a contradiction. Hence, the directional derivative (37) must be negative for all $\mathbf{x} \in \mathcal{Q}$, meaning that there exists a scalar $a > 0$ such that in the neighborhood $\mathcal{Q} \cap \mathcal{B}(\bar{\mathbf{x}}^{(n)}, a)$ the property (36) holds.

Combining (35) and (36), we obtain that, after a finite number k of iterations, the following condition holds:

$$\|\bar{\mathbf{x}}^{(n+k)} - \hat{\mathbf{x}}\| < \|\bar{\mathbf{x}}^{(n)} - \hat{\mathbf{x}}\|.$$

This implies that

$$\lim_{n \rightarrow \infty} \|\bar{\mathbf{x}}^{(n)} - \hat{\mathbf{x}}\| = 0,$$

and, for any $\sigma > 0$, we finally have

$$\exists n_\sigma : \min_{\mathbf{x}^{(i)} \in \mathbf{X}^{(n_\sigma)}} \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}\| < \sigma$$

thus proving the lemma. \square

Theorem 1. *Let Assumptions 1 and 3 hold. Then, $\forall \epsilon > 0, \exists n_\epsilon < \infty : z^{*(n_\epsilon)} \leq z^* + \epsilon$.*

Proof. Consider a point $\hat{\mathbf{x}} \in \mathcal{X}^*$ (see (2)) and take $\sigma = \frac{\epsilon}{\gamma_o}$. For any n , denote

$$\bar{\mathbf{x}}^{(n)} = \arg \min_{\mathbf{x}^{(i)} \in \mathbf{X}^{(n)}} \|\hat{\mathbf{x}} - \mathbf{x}^{(i)}\|$$

In virtue of Lemma 2, the exploration mode will be called infinitely often. Now, by applying Lemma 5, we have that $\exists n_\sigma < \infty : \|\hat{\mathbf{x}} - \bar{\mathbf{x}}^{(n_\sigma)}\| < \sigma = \frac{\epsilon}{\gamma_o}$.

Then, in virtue of Assumption 1 we have:

$$\begin{aligned}
f_o(\mathbf{x}^{*(n_\sigma)}) - f_o(\hat{\mathbf{x}}) &= z^{*(n)} - z^* \leq \\
f_o(\bar{\mathbf{x}}^{(n_\sigma)}) - f_o(\hat{\mathbf{x}}) &\leq \gamma_o \|\hat{\mathbf{x}} - \bar{\mathbf{x}}^{(n_\sigma)}\| < \epsilon \quad (38)
\end{aligned}$$

Thus proving the result with $n_\epsilon = n_\sigma$. \square

4.2. Optimality gap

Considering a run of the SMGO algorithm with N iterations, we now analyze the gap between the best sampled value and the actual global minimum,

$$\delta^{(N)} = z^{(N)*} - f_o(x^*)$$

denoted as the *optimality gap*. In calculating such a difference, the SM-based guarantees $\underline{z}^{(N)}(\mathbf{x}) \leq f_o(\mathbf{x}) \leq \bar{z}^{(N)}(\mathbf{x})$ can be utilized in the case of a known Lipschitz constant γ_o . Hence, the upper bound of the optimality gap $\bar{\delta}^{(N)}$ can be considered as the difference between $z^{(N)*}$ and the minimum lower bound $\underline{z}^{(N)}$ defined in (9) generated by data set $\mathbf{X}^{(N)}$, i.e.

$$\begin{aligned}
\bar{\delta}^{(N)} &= \max_{\mathbf{x} \in \mathcal{X}} \left(z^{(N)*} - \underline{z}^{(N)}(\mathbf{x}) \right) \\
&= z^{(N)*} - \min_{\mathbf{x} \in \mathcal{X}} \underline{z}^{(N)}(\mathbf{x}) \quad (39)
\end{aligned}$$

The calculation of $\min_{\mathbf{x} \in \mathcal{X}} \underline{z}^{(N)}(\mathbf{x})$ in (39) entails calculation of the hyperbolic Voronoi cells and of the lower bounds at their vertices. Denoting with $\mathbf{V}^{(N)}$ the set of such vertices generated from $\mathbf{X}^{(N)}$, (39) can be expressed as

$$\bar{\delta}^{(N)} = z^{(N)*} - \min_{\mathbf{v} \in \mathbf{V}^{(N)}} \underline{z}^{(N)}(\mathbf{v}). \quad (40)$$

The acquisition of HVD vertices using the samples can be approached as described in [24], which is referred for more information. The calculation of HVD cells and vertices entail computations with exponential complexity w.r.t. D [24]. However, an interested user can obtain an upper bound on the optimality gap by carrying out (40).

5. Computational Aspects

This section presents an analysis of the computational aspects of the proposed algorithm and possible improvements. The first aspect discussed is on ensuring the coverage of \mathcal{X} , followed by a discussion on the computational complexity.

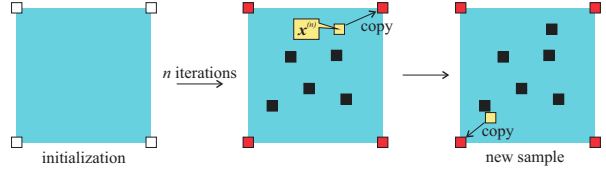


Figure 3: Corner mirroring of nearest sample point

5.1. On the coverage of the search set

The selection method for candidate points described in Section 3 allows SMGO to only explore the volume inside the convex hull of the sampled points in $\mathbf{X}^{(n)}$. Hence, in order to ensure exploration throughout all the search set \mathcal{X} , all its vertices $\mathbf{c}^{(v)}$ should also be considered in the candidate points generation, as stated in Assumption 3. One approach is to first evaluate f_o on all $\mathbf{c}^{(v)}$ before actually starting with exploitation/exploration. However, in the case of a hyperrectangle, this entails an exponential number of long function calls w.r.t. dimensionality D .

To address this issue, a simple approach is proposed in order to ensure coverage in \mathcal{X} without explicitly calling f_o to evaluate the corners. Moreover, with this approach Assumption 3 is not needed anymore for the theoretical results to hold. At each iteration n , the midpoints between any sample and all the vertices of \mathcal{X} are included as candidate exploration samples, even if the vertices do not belong to $\mathbf{X}^{(n)}$. Moreover, for the purpose of computing the uncertainty of each midpoint, the cost function value $f_o(\mathbf{c}^{(v)})$ at the v -th vertex $\mathbf{c}^{(v)}$ is estimated as:

$$f_o(\mathbf{c}^{(v)}) \approx z^{(w)}, \text{ where } w = \arg \min_{k=1, \dots, n} \|\mathbf{x}^{(k)} - \mathbf{c}^{(v)}\| \quad (41)$$

i.e., equal to the cost function value at the nearest sampled point. The concept of this strategy is shown in Fig. 3. Thanks to Assumption 1 and Lemma 5, with an increasing number of sample points in \mathcal{X} the estimate (41) approximates more and more accurately the true function value at the vertex.

5.2. On the computational burden

There are three main routines within an iteration of SMGO: Lipschitz constant update, exploitation by Mode θ routine, and eventually, if the expected improvement threshold is not achieved, exploration by Mode ψ . The worst-case SMGO complexity

analysis assumes that each iteration eventually proceeds to Mode ψ .

The first routine compares the incoming n -th sample with $n - 1$ existing ones, resulting in $\mathcal{O}(n)$ complexity.

The Modes θ and ψ can be designed as iterative routines as well to significantly improve the computational burden, as discussed in the following.

The exploitation Mode θ comprises the selection of $n - 1$ candidate points $\mathbf{x}_\theta^{(i)}$ for each segment $l_{n^*,i}$ from $\mathbf{x}^{*(n)}$ to all the other sampled points according to (15)-(16), and the calculation of $\underline{z}^{(n)}$ for each $\mathbf{x}_\theta^{(i)}$ using (9). Both are $\mathcal{O}(n)$ operations, leading to a compounded complexity of $\mathcal{O}(n^2)$. However, when $\gamma^{(n)}$ and $\mathbf{x}^{*(n)}$ do not change from iteration $n - 1$ to n , a cache can store $\mathbf{x}_\theta^{(i)}$ locations and corresponding $\underline{z}^{(n)}$. From (9), it is convenient to save $z^{(\tilde{k})}$ and $\tilde{\Delta}^{(i)} \triangleq \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(\tilde{k})}\|$, where

$$\tilde{k} = \arg \max_{k \in [1 \dots n]} \left(z^{(k)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(k)}\| \right).$$

In practice, information are saved¹ about the hypercone that actually generated $\underline{z}^{(n)}$ for each $\mathbf{x}_\theta^{(i)}$. As long as $\mathbf{x}^{*(n)}$ stays the same, previously existing $\mathbf{x}_\theta^{(i)}$ locations do not need to be computed again.

Furthermore, calculating $\underline{z}^{(n)}$ for an existing exploitation point $\mathbf{x}_\theta^{(i)}$ reduces to

$$\underline{z}^{(n)}(\mathbf{x}_\theta^{(i)}) = \max \left(\underline{z}^{(n-1)}(\mathbf{x}_\theta^{(i)}), \underline{z}_{new}^{(i)} \right) \quad (42)$$

where

$$\underline{z}_{new}^{(i)} = z^{(n)} - \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(n)}\|$$

for an incoming sample $(\mathbf{x}^{(n)}, z^{(n)})$.

Note that when $\underline{z}^{(n)}(\mathbf{x}_\theta^{(i)})$ changes due to (42), $z^{(\tilde{k})} = z^{(n)}$ and $\tilde{\Delta}^{(i)} = \mu\gamma^{(n)}\|\mathbf{x}_\theta^{(i)} - \mathbf{x}^{(n)}\|$. As a result, computing $\underline{z}^{(n)}$ for all existing candidate points is $\mathcal{O}(n)$, i.e. $\mathcal{O}(1)$ for each. Now, consider the situation when $\gamma^{(n)}$ is updated, assuming that $\mathbf{x}^{*(n)}$ and \tilde{k} do not change. This implies rescaling $\tilde{\Delta}^{(i)}$ such that

$$\tilde{\Delta}_{\langle n \rangle}^{(i)} = \frac{\gamma^{(n)}}{\gamma_{\langle n-1 \rangle}^{(n)}} \tilde{\Delta}_{\langle n-1 \rangle}^{(i)} \quad (43)$$

¹It can be understood that $z^{(\tilde{k})}$ is the value at the tip, while $\tilde{\Delta}^{(i)}$ is the height of the hypercone.

and updating $\underline{z}^{(n)}$ with the new sample $(\mathbf{x}^{(n)}, z^{(n)})$ using (42), still resulting in $\mathcal{O}(1)$ per candidate point. This implies a Mode θ aggregate complexity of $\mathcal{O}(n)$.

Finally, one new candidate point $\mathbf{x}_\theta^{(n)}$ is evaluated using (9), which is $\mathcal{O}(n)$. Hence, as long as $\mathbf{x}^{*(n)}$ does not change, the introduction of a cache reduces Mode θ to $\mathcal{O}(n)$. However, when $\mathbf{x}^{*(n)}$ changes, i.e. $\mathbf{x}^{*(n)} = \mathbf{x}^{(n)}$, computing all $\mathbf{x}_\theta^{(i)}$ locations must be repeated from the start, leading to $\mathcal{O}(n^2)$ worst-case complexity for Mode θ as stated before.

As defined, Mode ψ generates $\frac{n(n-1)}{2}$ candidate points (midpoints among all existing samples) at each iteration. Furthermore, the calculation of $\lambda^{(n)}$ for each candidate point is $\mathcal{O}(n)$, resulting in $\mathcal{O}(n^3)$ for each exploration iteration.

Note that behavior in Mode ψ does not depend on the current best point, because n new segments are added to the search set at iteration n . A cache solution can be introduced also in this mode, updating $\lambda^{(n)}$ -related values iteratively.

For each existing $\mathbf{x}_\psi^{(i,j)}$, the values to be saved are $z^{(\hat{k})}$ and $\hat{\Delta}^{(i,j)}$ (for $\bar{z}^{(n)}$), and $z^{(\check{k})}$ and $\check{\Delta}^{(i,j)}$ (for $\underline{z}^{(n)}$), where \hat{k} and \check{k} are defined as

$$\hat{k} = \arg \min_{k \in [1 \dots n]} \left(z^{(k)} + \mu\gamma^{(n)}\|\mathbf{x}_\psi^{(i,j)} - \mathbf{x}^{(k)}\| \right),$$

$$\check{k} = \arg \max_{k \in [1 \dots n]} \left(z^{(k)} - \mu\gamma^{(n)}\|\mathbf{x}_\psi^{(i,j)} - \mathbf{x}^{(k)}\| \right).$$

The iterative updates are performed using eq. (42) and (43). This then results in $\mathcal{O}(1)$ for each existing midpoint, and $\mathcal{O}(n)$ for each new midpoint generated by new segments from incoming $\mathbf{x}^{(n)}$ to the $n - 1$ existing points. This results in $\mathcal{O}(n^2)$ for Mode ψ , and in turn, a $\mathcal{O}(n^2)$ worst-case complexity for SMGO.

Remark 1. Note that when $\gamma^{(n)}$ is updated (which always implies a growth of its value), the proposed iterative implementation might result in more conservative bounds $\bar{z}^{(n)}$ and $\underline{z}^{(n)}$ than if repeatedly recalculated at every iteration, i.e. lower $\underline{z}^{(n)}$ for Mode θ , and larger $\lambda^{(n)}$ in Mode ψ . This occurs because with higher $\gamma^{(n)}$ -values the HVD layout increasingly approaches that of a non-hyperbolic Voronoi tessellation, i.e. the cones providing the tightest bounds at any $\mathbf{x}_\psi^{(i,j)}$ could change when

$\gamma^{(n+1)} \gg \gamma^{(n)}$, and $z^{(\hat{k})}$ and $z^{(\tilde{k})}$ would tend to the nearest sample w.r.t. said $\mathbf{x}_\psi^{(i,j)}$. These effects are not accounted for by the iterative update procedure, which assumes that the bound-generating cones remain the same. However, the theoretical properties of SMGO are still valid with the iterative implementation. In fact, all the arguments for Lemma 2 are valid also with iteratively-computed lower bounds. Regarding the exploration mode, inequality (26) in Lemma 4 is consistent with the more conservative iteratively-computed $\gamma^{(n)}$. Furthermore, when sampling a new point $\mathbf{m}_1^{(n)}$, for which $\lambda^{(n+1)}(\mathbf{m}_1^{(n)})$ is exactly computed, (33) still holds. Hence, all arguments laid out in Lemma 5 still hold.

6. Performance Test Results

In this section, the proposed SMGO algorithm is evaluated on well-known benchmarks, usually employed in black-box optimization, comparing its performance with representative Lipschitz-based methods: DIRECT and AdaLIPO. In addition, the Bayesian optimization approach (BayesianO) is also considered in the tests. The results are discussed in terms of iteration-based optimization performance and computational times. All scripts used for this section are available at <https://github.com/lorenzozabugjr/smgo>.

6.1. Test parameters

Seven test functions of varying structure and dimensions are chosen for this test, covering a variety of characteristics relevant to black-box optimization benchmarks. Their search bounds, optimal values and relevant features are summarized in Table 1. A reasonable assumption in real-life applications is that the black-box function has no available optimal value. Hence, each algorithm is given a budget of $N = 500$ long function evaluations for each optimization run. The optimal results after the evaluation budget are then measured and compared among the different algorithms, in what is referred as a fixed-cost/budget comparison [25]. All computations were performed in MATLAB 2020b on a system with AMD Ryzen 9 3900X (3.80 GHz) and 32 GB RAM.

To assess the sensitivity of each algorithm to different initial information and randomized decisions, 100 independent runs (trials) were performed. For each trial, the same randomly-generated starting point is given to AdaLIPO, SMGO, and BayesianO.

On the other hand, DIRECT considers a fixed set of search points, not allowing to randomize the starting samples.

SMGO is implemented as described in Algorithm 1 with parameters $\alpha = 0.001$ and $\mu = 1.025$. AdaLIPO is implemented as described in [18] with parameter $p = 0.1$. DIRECT is implemented as described in [14], not requiring any tuning parameter. Finally, the algorithm available in the Matlab Statistics and Machine Learning Toolbox is employed for the BayesianO implementation, using the provided default parameters.

6.2. Iteration-based performance

The evolution of the best sample $z^{*(n)}$ along the iterations ($n \in [50, 500]$) for each algorithm is shown in Figs. 4 and 5 for the Rosenbrock and Deb’s functions, respectively. The graphs show the distribution of the best sampled cost for the 100 independent runs. Note that the results of the DIRECT algorithm do not show dispersion due to the fact that it is a deterministic search method. DIRECT follows a batch sampling strategy, hence the graphs show for each n value, the best results for the previously-sampled batches.

The Rosenbrock function is a standard test case for gradient-based algorithms, with a unique global minimum (unimodal) for $D = 2$; however it has 2 minima for $D = 4 \sim 30$ [26]. The location of \mathbf{x}^* is at $(1, 1, \dots, 1)$, hence the classical search bounds $-10 \leq x_i \leq 10$ would be advantageous for DIRECT, which invariably samples the center of the search hyper-box. For this reason, the search bounds are changed to move \mathbf{x}^* away from the center.

As seen in Fig. 4, DIRECT has only improved marginally the initial sampled cost after the 500 iterations. Also AdaLIPO achieved marginal improvements after 500 iterations. This might be associated to the completely randomized exploration and exploitation decisions. On the other hand, both SMGO and BayesianO achieved large improvements during the first 100 iterations, while slight enhancements are observed during the final 400 iterations. SMGO exhibits a wider distribution among different runs and, on average, achieves a better result after 500 iterations.

In the case of Deb’s function #1, as shown in Fig. 5, the presence of numerous global minima causes a bad performance for DIRECT, which barely improved its initial best cost. On the other

Description	Function definition $f(\mathbf{x})$	Bounds	z^*	Features
Rosenbrock	$\sum_{i=1}^D [100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2]$	$-40 \leq x_i \leq 5$	0.0	unimodal, non-separable
Styblinski-Tang	$\frac{1}{2} \sum_{i=1}^D (x_i^4 - 16x_i^2 + 5x_i)$	$-5 \leq x_i \leq 5$	$-39.166D$	multimodal, separable
Deb's #1	$\frac{1}{D} \sum_{i=1}^D \sin^6(5\pi x_i)$	$-1 \leq x_i \leq 1$	-1.0	multimodal, separable, numerous global minima
Deb's #2	$\frac{1}{D} \sum_{i=1}^D \sin^6[5\pi(x_i^{3/4} - 0.05)]$	$0 \leq x_i \leq 150$	-1.0	multimodal, separable, numerous global minima
Schwefel	$-\sum_{i=1}^D x_i \sin(\sqrt{ x_i })$	$-500 \leq x_i \leq 500$	$-418.982D$	multimodal, separable, numerous local minima
Salomon	$1 - \cos\left(2\pi\sqrt{\sum_{i=1}^D x_i^2}\right) + 0.1\sqrt{\sum_{i=1}^D x_i^2}$	$-40 \leq x_i \leq 70$	0	multimodal, non-separable
Brown	$\sum_{i=1}^{D-1} (x_i^2)^{(x_{i+1}^2+1)} + (x_{i+1}^2)^{(x_i^2+1)}$	$-1 \leq x_i \leq 4$	0	multimodal, non-separable

Table 1: Test functions used for comparative tests

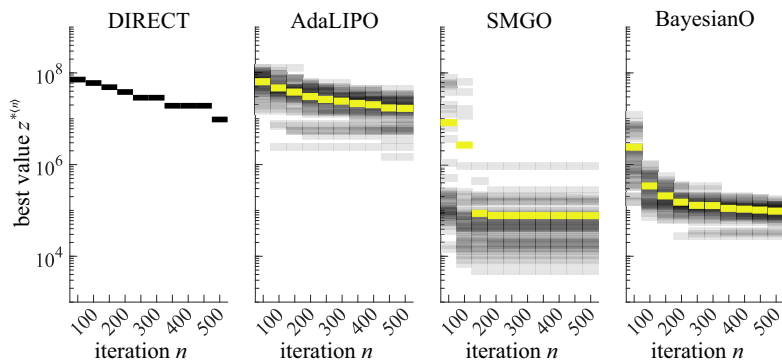


Figure 4: Optimal value distribution w.r.t. iterations, 10D Rosenbrock function (log scale)

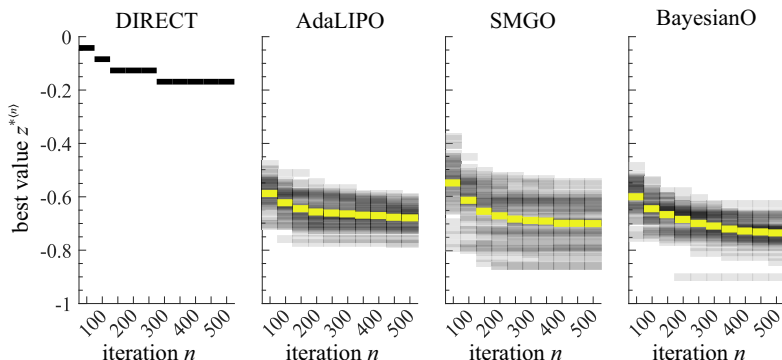


Figure 5: Optimal value distribution w.r.t. iterations, 10D Deb's function #1

hand, AdaLIPO, BayesianO, and SMGO have resulted in similar results, achieving very similar average costs after the iteration limit, while BayesianO shows more concentrated values among runs. In this case the randomized mechanism of AdaLIPO allows it to discover low-cost regions to improve the best sample. Furthermore, Bayesian’s usage of a response surface, and SMGO’s geometry-based candidate points derivation have become useful in looking for low-cost regions in the search space.

A summary (mean best values) of the algorithms’ outcomes after 500 iterations are shown in Table 2. Note that the results tabulated for DIRECT were achieved at iterations which are a bit larger than the pre-defined maximum $N = 500$, owing to its batch-based samples generation. In most test cases, BayesianO achieved the best average results, while SMGO and DIRECT attained the best average performance in 2 cases each. On the other hand, AdaLIPO achieved fair results, even if it ranked second-place or worse for most of the tested functions.

Furthermore, for each test function (with the same starting points across all algorithms), the number of trials (out of 100) resulting in best results for each respective algorithm is shown in Table 3. As with the average performance, BayesianO achieved the most best trials in 9 out of the 13 tested functions. Furthermore, in Deb’s #1 (5D), Schwefel (5D and 10D), and Salomon (5D and 10D), it resulted in best results in almost all of the given trials. On the other hand, SMGO had the most best results in Rosenbrock and Deb’s #2 (10D), and had the second largest count of best trials in Deb’s #1 (10D) and Deb’s #2 (5D). AdaLIPO registered the second count of best trials in Deb’s #2 (10D), and had no best trials for 4 functions. Lastly, DIRECT had no best trials for 8 out of 13 test functions, but took almost all the best results in the Brown function (5D and 10D).

We have also compared the outcomes of SMGO with the competitor methods using the Wilcoxon and Kruskal-Wallis non-parametric statistical tests, pairing the results by the shared random initial point. Both test were performed with 5% significance level and with the null hypothesis that SMGO results are statistically similar to those from the respective competitor method. The two tests agreed that SMGO results were statistically similar to AdaLIPO for Styblinski-Tang (5D), Deb’s #2 (5D, 10D), Schwefel (10D), and Brown (10D). Furthermore, similarities were found between SMGO

and BayesianO for Deb’s #2 (10D). Given these results, SMGO is found to have competitive results compared with the considered global optimization algorithms on the test functions employed in the analysis.

6.3. Computational time

The computational time required to complete an iteration for each considered algorithms is analyzed for the case of the 5D Deb’s function #1 optimization. Fig. 6 (in log scale) shows the time required to complete an iteration as n increases. The mechanism of DIRECT calculates recursive hyperboxes to be sampled by batches, hence the computational time *per iteration* is not well defined, therefore results are not shown in the plot. For AdaLIPO an almost flat line is observed, with 10 to 100 μ s per iteration, implying a computational time independent of the iteration number. SMGO has resulted in a polynomial complexity w.r.t. iterations, which is expected from the computational analysis presented in Section 5. Finally, BayesianO is the most demanding algorithm, requiring computational times per iteration between 2 and 3 orders of magnitude larger than SMGO in this test case.

The computational times observed for all the considered test functions have shown trends similar to those on the 5D Debs function #1. Table 4 summarizes the results, showing the total computational time required to complete an optimization run. The AdaLIPO algorithm, based on a simple mechanism of randomized generation of exploitation and exploration points, always shows very fast computational times, not affected by the dimensionality of the decision variable, scarcely affected by n , and taking 0.01-0.02 s per optimization run for all the test functions.

From the tests, BayesianO requires the largest computational times, which achieved around 4-5 s per iteration at $n = 500$ for most of the 10 dimensions functions. On the other hand, SMGO exhibits an intermediate per-iteration time around 100 μ s to 10 ms for 5 dimensions functions and tens to hundreds of ms for 10 dimensions functions. Furthermore, for SMGO the total optimization run time when applied to 10 dimensions functions is around 10 times those employed for 5 dimensions functions. However, this is significantly faster than BayesianO, ranging from around 15 times shorter time for 10D functions, to 70 times shorter time for 5D functions.

The observed computational times, combined with the competitive optimization results presented

Test function	D	DIRECT	AdaLIPO	SMGO	BayesianO
Rosenbrock	10	1.17 E+5	1.77 E+7	8.63 E+4	9.35 E+4
Styblinski-Tang	5	-1.95 E+2	-1.59 E+2	-1.58 E+2	-1.95 E+2
	10	-3.19 E+2	-2.61 E+2	-2.96 E+2	-3.33 E+2
Deb's #1	5	-5.29 E-1	-8.29 E-1	-8.07 E-1	-9.68 E-1
	10	-2.11 E-1	-6.74 E-1	-6.97 E-1	-7.35 E-1
Deb's #2	5	-5.97 E-1	-8.32 E-1	-8.33 E-1	-8.59 E-1
	10	-3.98 E-1	-6.72 E-1	-6.81 E-1	-6.77 E-1
Schwefel	5	-1.47 E+3	-1.28 E+3	-1.23 E+3	-1.90 E+3
	10	-1.48 E+3	-1.83 E+3	-1.79 E+3	-2.67 E+3
Salomon	5	3.45 E0	2.64 E0	2.19 E0	6.38 E-1
	10	5.05 E0	5.79 E0	5.29 E0	2.40 E0
Brown	5	6.46 E-4	1.39 E-1	8.29 E-2	3.90 E-3
	10	4.76 E-2	9.78 E-1	9.61 E-1	1.65 E-1

Table 2: Results summary for comparative tests: averages of $z^{*(n)}$ over 100 runs, after 500 iterations per trial.

Test function	D	DIRECT		AdaLIPO		SMGO		BayesianO	
		1st	2nd	1st	2nd	1st	2nd	1st	2nd
Rosenbrock	10	3	34	0	0	75	12	22	54
Styblinski-Tang	5	12	88	0	0	0	0	88	12
	10	15	71	0	0	8	15	77	14
Deb's #1	5	0	0	2	61	2	35	96	4
	10	0	0	12	37	37	25	51	38
Deb's #2	5	0	0	27	28	28	36	45	36
	10	0	0	29	28	44	25	27	47
Schwefel	5	0	82	1	8	0	9	99	1
	10	0	0	0	56	0	44	100	0
Salomon	5	0	0	0	30	0	70	100	0
	10	0	48	0	8	0	44	100	0
Brown	5	98	2	0	0	0	0	2	98
	10	100	0	0	0	0	0	0	100

Table 3: Results summary for comparative tests: number of best and 2nd best trials for the tested algorithms

Test function	D	AdaLIPO	SMGO	BayesianO
Rosenbrock	10	0.017	74.96	1071.04
Styblinski-Tang	5	0.013	7.71	511.66
	10	0.014	75.01	1085.11
Deb's #1	5	0.012	7.71	574.08
	10	0.014	74.64	1164.80
Deb's #2	5	0.012	7.72	576.63
	10	0.014	73.97	1176.10
Schwefel	5	0.011	7.77	583.24
	10	0.014	75.13	1188.08
Salomon	5	0.013	7.63	643.25
	10	0.015	75.08	1033.54
Brown	5	0.012	7.82	476.40
	10	0.013	75.47	886.81

Table 4: Comparison of total optimization run times, averaged over the 100 independent runs for each test case (all in seconds)

in Figs. 4-5, and Table 2, indicate that the proposed SMGO algorithm shows a good trade-off between optimization performance and computational speed, compared to the other considered methods.

7. Conclusions and Further Work

In this work a sequential algorithm for global optimization of black-box functions has been proposed. The SMGO algorithm assumes a Lipschitz-

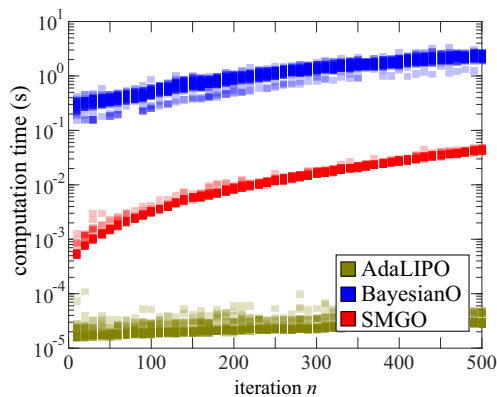


Figure 6: Calculation times on 5D Deb's #1 (log scale)

continuous cost function and is based on a nonlinear Set Membership function approximation. The selection of the test points where the black-box function is evaluated is carried out by solving simplified geometry-based problems on the guaranteed lower bounds or uncertainty intervals of the function, provided by a Set Membership model. It is shown that the SMGO algorithm finds a sub-optimal solution with desired tolerance after finite iterations and converges asymptotically to the global minimum. Furthermore, the computational complexity of the algorithm is discussed and a cache-based solution is introduced to improve the computational burden.

The proposed method is evaluated against other representative global optimization methods on several test functions, comparing the quality of the best solutions found after a fixed number of calls to the black-box function. The results show the competitiveness of the SMGO algorithm compared to state of the art methods found in literature. Ongoing research aims to extend the algorithm to include constraints and its evaluation on experimental applications.

References

- [1] D. R. Jones, M. Schonlau, and W. J. Welch, "Efficient Global Optimization of Expensive Black-Box Functions," *Journal of Global Optimization*, vol. 13, pp. 455–492, 1998.
- [2] N. Pham, A. Malinowski, and T. Bartczak, "Comparative study of derivative free optimization algorithms," *IEEE Transactions on Industrial Informatics*, vol. 7, no. 4, pp. 592–600, 2011.
- [3] S. Amaran, N. V. Sahinidis, B. Sharda, and S. J. Bury, "Simulation optimization: a review of algorithms and applications," *Annals of Operations Research*, vol. 240, no. 1, pp. 351–380, 2016.
- [4] S. Gao, L. Shi, and Z. Zhang, "A peak-over-threshold search method for global optimization," *Automatica*, vol. 89, pp. 83 – 91, 2018.
- [5] J. Kennedy and R. Eberhart, "Particle swarm optimization," in *Proceedings of ICNN'95 - International Conference on Neural Networks*, vol. 4. IEEE, 2014, pp. 1942–1948.
- [6] Q. Liu, S. Yang, and J. Wang, "A collective neurodynamic approach to distributed constrained optimization," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 8, pp. 1747–1758, 2017.
- [7] H. Han, W. Lu, and J. Qiao, "An Adaptive Multiobjective Particle Swarm Optimization Based on Multiple Adaptive Methods," *IEEE Transactions on Cybernetics*, vol. 47, no. 9, pp. 2754–2767, 2017.
- [8] M. J. Powell, "UOBYQA: Unconstrained optimization by quadratic approximation," *Mathematical Programming, Series B*, vol. 92, no. 3, pp. 555–582, 2002.
- [9] A. Bemporad, "Global optimization via inverse distance weighting and radial basis functions," *Computational Optimization and Applications*, vol. 77, no. 2, pp. 571–595, 2020.
- [10] Y. Wang, D. Q. Yin, S. Yang, and G. Sun, "Global and local surrogate-assisted differential evolution for expensive constrained optimization problems with inequality constraints," *IEEE Transactions on Cybernetics*, vol. 49, no. 5, pp. 1642–1656, 2019.
- [11] S. Piyavskii, "An algorithm for finding the absolute extremum of a function," *USSR Computational Mathematics and Mathematical Physics*, vol. 12, no. 4, pp. 57 – 67, 1972.
- [12] B. O. Shubert, "A sequential method seeking the global maximum of a function," *SIAM Journal on Numerical Analysis*, vol. 9, no. 3, pp. 379–388, 1972.
- [13] D. R. Jones, C. D. Perttunen, and B. E. Stuckman, "Lipschitzian optimization without the Lipschitz constant," *Journal of Optimization Theory and Applications*, vol. 79, no. 1, pp. 157–181, oct 1993.
- [14] D. E. Finkel and C. T. Kelley, "Additive Scaling and the DIRECT Algorithm," *Journal of Global Optimization*, vol. 36, no. 4, pp. 597–608, oct 2006.
- [15] S. Z. Khong, D. Nešić, C. Manzie, and Y. Tan, "Multi-dimensional global extremum seeking via the direct optimization algorithm," *Automatica*, vol. 49, no. 7, pp. 1970 – 1978, 2013.
- [16] R. Paulavičius and J. Žilinskas, "Simplicial Lipschitz optimization without the Lipschitz constant," *Journal of Global Optimization*, vol. 59, no. 1, pp. 23–40, may 2014.
- [17] R. Paulavičius, Y. D. Sergeyev, D. E. Kvasov, and J. Žilinskas, "Globally-biased DISIMPL algorithm for expensive global optimization," pp. 545–567, 2014.
- [18] C. Malherbe and N. Vayatis, "Global optimization of Lipschitz functions," *34th International Conference on Machine Learning, ICML 2017*, vol. 5, pp. 3592–3601, 2017.
- [19] M. Milanese and C. Novara, "Unified set membership theory for identification, prediction and filtering of nonlinear systems," *Automatica*, vol. 47, no. 10, pp. 2141 – 2151, 2011.
- [20] —, "Set Membership identification of nonlinear systems," *Automatica*, vol. 40, no. 6, pp. 957–975, 2004.
- [21] C. Novara, F. Ruiz, and M. Milanese, "Direct filtering: A new approach to optimal filter design for nonlinear systems," *IEEE Transactions on Automatic Control*,

- vol. 58, no. 1, pp. 86–99, 2013.
- [22] L. Fagiano and C. Novara, “Learning a Nonlinear Controller From Data: Theory, Computation, and Experimental Results,” *IEEE Transactions on Automatic Control*, vol. 61, no. 7, pp. 1854–1868, jul 2016.
 - [23] L. Sabug, F. Ruiz, and L. Fagiano, “On the use of Set Membership theory for global optimization of black-box functions,” in *Proceedings of the 59th Conf. Decision and Control (CDC)*. IEEE, 2020.
 - [24] M. Milanese and C. Novara, “Computation of local radius of information in SM-IBC identification of nonlinear systems,” *Journal of Complexity*, vol. 23, no. 4-6, pp. 937–951, 2007.
 - [25] V. Beiranvand, W. Hare, and Y. Lucet, “Best practices for comparing optimization algorithms,” *Optimization and Engineering*, vol. 18, no. 4, pp. 815–848, sep 2017.
 - [26] Y.-W. Shang and Y.-H. Qiu, “A Note on the Extended Rosenbrock Function,” *Evolutionary Computation*, vol. 14, no. 1, pp. 119–126, mar 2006.