# META-SMGO-$\triangle$: similarity as a prior in black-box optimization

Riccardo Busetto[1], Valentina Breschi[2] and Simone Formentin[1]

*Abstract*—**When solving global optimization problems in practice, one often ends up repeatedly solving problems that are similar to each others. By introducing a rigorous definition of similarity to exploit priors obtained from past experience to efficiently solve new (similar) problems, in this work we incorporate the META-learning rationale into SMGO-$\triangle$, a global optimization approach recently proposed in the literature. Through a benchmark numerical example we show the practical benefits of our META-extension of the baseline algorithm, while providing theoretical bounds on its performance.**

## I. INTRODUCTION

Black-box optimization is a fundamental tool whenever objective functions and/or constraints are unknown or expensive to evaluate. Such an approach to optimization has been studied extensively, with many different algorithms proposed for solving this problem. As a noticeable example, Bayesian Optimization uses probabilistic models to guide the search process and has proven its effectiveness in many applications [1]. Other methods include gradient-based optimization [2], simulated annealing [3], and particle swarm optimization [4].

Despite the large number of available algorithms, black-box optimization still remains a challenging problem, particularly in high-dimensional spaces or when the objective function is noisy or non-convex. Specifically, when the vector of parameters is large, one might require several iterations before converging to a satisfactory solution. To avoid this, practitioners usually acquire expertise by repeatedly solving problems that despite their differences share some common features and, thus, they are somehow *similar*. The idea behind this paper starts from the observation that the experience acquired in solving optimization problems with shared characteristics can be re-used as valuable prior to solve more efficiently a new (but similar) problem. This intuition is also at the foundation of *meta-learning*, a sub-field of machine learning that focuses on developing algorithms that can automatically learn how to solve new tasks more efficiently and effectively by leveraging prior experience from similar tasks [5], [6].

There are several approaches to meta-learning [7] including metric-based, model-based, and optimization-based methods. Nonetheless, while meta-learning has a wide range of applications in machine learning, see e.g., [8], [9], [10],

[1]R. Busetto and S. Formentin are with Dip. di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Milano, Italy. `name.surname@polimi.it`

[2]V. Breschi is with Dept. of Electrical Engineering, Eindhoven University of Technology, Eindhoven, The Netherlands . `v.breschi@tue.nl`

it has seldom been exploited for control systems, with just a few very recent exceptions [11], [12], [13], [14], despite the fact that it could of great benefit in several control applications. Indeed, many control algorithms, especially those based on the use of experimental data, require the calibration of hyper-parameters that significantly affect the resulting closed-loop performance [15], which are generally chosen with expensive experiments, involving either sensitivity analyses, or driven by expert-based and rule-of-thumb-based design.

Because of its remarkable efficiency when compared to similar approaches like Bayesian optimization, in this work we consider a recently proposed black-box optimization technique called SMGO-$\triangle$ [16] and we apply meta-learning tools to show that prior experience with similar problems can be used to further boost its performance. In particular, we leverage on the set-membership nature of SMGO-$\triangle$ to derive a theoretical bound on the performance of its META-version, and we show on a numerical case study that the META-version of SMGO-$\triangle$ reduces both the number of iterations for convergence and constraint violations.

The remainder of the paper is organized as follows. For the self-consistency of the work, Section II is devoted to a review of the main features of SMGO-$\triangle$. These preliminaries allows us to formulate the problem stated in Section III and to discuss the proposed META extension of SMGO-$\triangle$, introduced in Section IV. In this section, we further prove two key properties of META-SMGO-$\triangle$, while discussing how the proposed approach can be implemented in practice. The effectiveness of META-SMGO-$\triangle$ is then evaluated in Section V, showing the advantages of exploiting similarities for hyper-parameter tuning for SMGO-$\triangle$. The paper is ended by some concluding remarks.

## II. AN OVERVIEW ON SMGO-$\triangle$

SMGO-$\triangle$ is an iterative optimization procedure specifically devised to tackle the following class of problems

$$\min_{X \in \mathcal{X}} \quad f(X) \tag{1a}$$

$$\text{s.t. } g_s(X) \geq 0, \quad s = 1, \dots, S, \tag{1b}$$

where the function $f : \mathcal{X} \mapsto \mathbb{R}$ one aims at minimizing is assumed to be *unknown*. Note that the minimum, namely

$$X^\star = \underset{X \in \mathcal{X}, g_s(X) \geq 0}{\arg\min} f(X), \tag{2}$$

should satisfy a set of inequality constraints, characterized by $g_s(X) : \mathcal{G}_s \mapsto \mathbb{R}$, for $s = 1, \dots, S$, that are also supposed to be *unknown*. The method rests on the following assumptions.

*Assumption 1:* The functions $f(\cdot)$ and $g_s(\cdot)$ are assumed to be Lipschitz continuous, namely they satisfy:

$$|h(X_1) - h(X_2)| \leq \gamma_h \|X_1 - X_2\|_2, \quad \forall X_1, X_2 \in \mathcal{X}, \quad (3)$$

with $h(\cdot)$ being a placeholder for either of the two functions, and $\gamma_h > 0$ being the associated Lipschitz constant.

*Assumption 2:* The search space $\mathcal{X}$ and the space of all feasible solutions are not disjoint, namely

$$\mathcal{X} \cap \left\{ \bigcap_{s=1}^{S} \mathcal{G}_s \right\} \neq \emptyset, \quad (4)$$

thus implying that a solution to (1) exists.

By relying on these assumptions, SMGO-$\Delta$ addresses the problem in (1) by iteratively performing three following steps, until a maximum number of iterations $n_{\max}$ is attained.

Let $\mathcal{D}^{(n)}$ be the set comprising all the information on the points $X$ explored up to the $n-1$-th iteration, i.e., the points and the corresponding values of $f(\cdot)$ and $g_s(\cdot)$ defined as:

$$z^{(n)} = f\left(X^{(n)}\right), \quad (5)$$

$$c_s^{(n)} = g_s\left(X^{(n)}\right). \quad (6)$$

The information in $\mathcal{D}^{(n)}$ are initially used in SMGO-$\Delta$ to estimate the Lipschitz constant of the objective function as

$$\gamma_f^{(n)} = \max_{\substack{(X^{(i)}, z^{(i)}) \in \mathcal{D}^{(n)}, \\ (X^{(j)}, z^{(j)}) \in \mathcal{D}^{(n)}}} \left( \frac{|z^{(i)} - z^{(j)}|}{\|X^{(i)} - X^{(j)}\|}, \underline{\gamma}_f \right), \quad (7)$$

and that of each function characterizing the constraints as

$$\gamma_{g_s}^{(n)} = \max_{\substack{(X^{(i)}, c_s^{(i)}) \in \mathcal{D}^{(n)}, \\ (X^{(j)}, c_s^{(j)}) \in \mathcal{D}^{(n)}}} \left( \frac{|c_s^{(i)} - c_s^{(j)}|}{\|X^{(i)} - X^{(j)}\|}, \underline{\gamma}_{g_s} \right), \quad (8)$$

for $s = 1, \ldots, S$. Note that both computation rely on a lower-bound on the Lipschitz constant (i.e., $\underline{\gamma}_f$ and $\{\underline{\gamma}_{g_s}\}_{s=1}^{S}$), initialized by the user and then iteratively replaced with the updated estimate of these constants over exploration[1]. These estimates are used to update the *bounding functions*

$$\overline{f}^{(n)}(X) = \min_{k=1,\ldots,n} \left( z^{(k)} + \gamma_f^{(n)} \left\| X - X^{(k)} \right\| \right), \quad (9a)$$

$$\underline{f}^{(n)}(X) = \max_{k=1,\ldots,n} \left( z^{(k)} - \gamma_f^{(n)} \left\| X - X^{(k)} \right\| \right), \quad (9b)$$

$$\overline{g}_s^{(n)}(X) = \min_{k=1,\ldots,n} \left( c_s^{(k)} + \gamma_{g_s}^{(n)} \left\| X - X^{(k)} \right\| \right), \quad (9c)$$

$$\underline{g}_s^{(n)}(X) = \max_{k=1,\ldots,n} \left( c_s^{(k)} - \gamma_{g_s}^{(n)} \left\| X - X^{(k)} \right\| \right), \quad (9d)$$

for each $s \in \{1, \ldots, S\}$, the *central approximations*

$$\tilde{f}^{(n)}(X) = \frac{1}{2} \left( \overline{f}^{(n)}(X) + \underline{f}^{(n)}(X) \right), \quad (10a)$$

$$\tilde{g}_s^{(n)}(X) = \frac{1}{2} \left( \overline{g}_s^{(n)}(X) + \underline{g}_s^{(n)}(X) \right), \quad s = 1, \ldots, S. \quad (10b)$$

and the *uncertainty bounds*

$$\lambda_f^{(n)}(X) = \overline{f}^{(n)}(X) - \underline{f}^{(n)}(X), \quad (11a)$$

$$\lambda_{g_s}^{(n)}(X) = \overline{g}_s^{(n)}(X) - \underline{g}_s^{(n)}(X), \quad s = 1, \ldots, S. \quad (11b)$$

The central approximation and the uncertainty bounds are then exploited to choose the next candidate point $X_\theta^{(n)}$, which is obtained by solving

$$\min_{X \in E^{(n)} \cup \mathcal{T}^{(n)}} \quad \tilde{f}^{(n)}(X) - \beta \lambda_f^{(n)}(X) \quad (12a)$$

$$\text{s.t.} \quad \Delta \tilde{g}_s^{(n)}(X) + (1-\Delta)\underline{g}_s^{(n)}(X) \geq 0, \ \forall s, \quad (12b)$$

where $\beta, \Delta > 0$ are tunable parameters, and the state of space is restricted to the intersection of a (cumulative) set of candidate points $E^{(n)}$ and a *trust region* $\mathcal{T}^{(n)}$ centered on the current estimate $X^{\star(n)}$ of the minima (2), defined as

$$\mathcal{T}^{(n)} = \left\{ X \in \mathcal{X} : \left\| X - X^{\star(n)} \right\|_2 \leq v^{(n)} \right\}. \quad (13)$$

An *expected improvement* test is then performed on $X_\theta^{(n)}$, by checking if the bounding function $\underline{f}^{(n)}(\cdot)$ satisfies the following:

$$\underline{f}^{(n)}(X_\theta^{(n)}) \leq z^{\star(n)} - \alpha \gamma_f^{(n)}, \quad (14)$$

where $z^{\star(n)}$ is the function value at the current estimated minima, and $\alpha > 0$ is another tunable parameter. If the inequality is satisfied, then $X_\theta^{(n)}$ becomes the point to be evaluated at the next iteration

$$X^{(n+1)} \leftarrow X_\theta^{(n)}, \quad (15)$$

otherwise *exploration* is promoted over more uncertain regions by defining a new candidate point $X_\psi^{(n)}$ as

$$X_\psi^{(n)} = \arg\max_{X \in E^{(n)}} \ \xi_\psi^{(n)}(X), \quad (16)$$

where $\xi_\psi^{(n)} : E^{(n)} \to \mathbb{R}$ is the so-called *exploration merit function*. Since this specific step is not modified by our META extension of the approach, we refer the reader to [16] for additional details.

## III. PROBLEM STATEMENT

Let us assume that we have already exploited SMGO-$\Delta$ to retrieve the global optimum of $M$ functions, that are *similar* to the one we aim at optimizing. Our goal is to exploit this similarity to enhance the performance of SMGO-$\Delta$ in tackling this new optimization instance. To formally state this problem, let us define the concept of *similarity* considered in this work.

*Definition 1:* Let $f_1(X, \varphi_1)$ and $f_2(X, \varphi_2)$ be two functions satisfying Assumption 1. These functions are said:

- $\rho$-similar, if there exists $\rho < \infty$ such that the radius $\overline{\rho}$ of the *smallest enclosing circle* of their minima $X_1^\star$ and $X_2^\star$ satisfies $\overline{\rho} \leq \rho$.
- $\zeta$-similar, if there exists $\zeta \in [0, \infty)$ such that $\overline{\zeta} \leq \zeta$, with $\overline{\zeta}$ being the distance of their Lipschitz constants:

$$\overline{\zeta} = |\gamma_1 - \gamma_2|. \quad (17)$$

---

[1] At the first iteration of SMGO-$\Delta$, the estimates of the Lipschitz constants are set at a user-defined lower-bound.

Suppose now that $M$ constrained problems in the form of (1) are solved with SMGO-$\Delta$ over the same number of optimization steps $n_{max}$, for different instances of the cost function that are yet $\rho$-similar and $\zeta$-similar according to Definition 1. Let us additionally assume that all these problems share the same set of constraints[2]. As outcomes of these optimization routines, one can extract the resulting estimated minima $X_i^\star$ and minimum value $z_i^\star$, the set of explored states $\{X_i^{(n)}\}_{n=1}^{n_{max}}$, function values $\{z_i^{(n)}\}_{n=1}^{n_{max}}$ and constraints $\{c_{s,i}^{(n)}\}_{n=1}^{n_{max}}$, for $s = 1, \ldots, S$, and the estimates of the Lipschitz contants

$$\hat{\gamma}_i = \begin{bmatrix} \gamma_{f,i}^{(n_{max})} & \gamma_{g_1,i}^{(n_{max})} & \cdots & \gamma_{g_S,i}^{(n_{max})} \end{bmatrix}^\top, \quad (18)$$

with $i = 1, \ldots, M$. These elements can be used to construct a *META-dataset* $\mathcal{D}^{\mathrm{meta}}$, that, in turn, can be exploited to tackle a new optimization problem (1) when the cost function $f(X)$ is $\rho$-similar and $\zeta$-similar to the $M$ functions already optimized.

Under our assumptions, we thus aim at exploiting $\mathcal{D}^{\mathrm{meta}}$ to $(i)$ reduce the number of iterations required by SMGO-$\Delta$ to find the global optimum, and $(ii)$ reduce the number of constraints violations throughout the optimization. In this work, this goal is achieved by using $\mathcal{D}^{\mathrm{meta}}$ to initialize both the first evaluation point $X^{(1)}$ of the new instance of SMGO-$\Delta$ and the initial lower bounds on the Lipschitz constant $\underline{\gamma}_f$ required at the first to compute (9)-(10), i.e., to **META-initialize** the new instance of SMGO-$\Delta$.

## IV. META-LEARNING FOR SMGO-$\Delta$

The performance of SMGO-$\Delta$ are, by construction, shaped by the initial choices that the user has to perform. Here we focus on two crucial hyper-parameters, namely the initial exploration point and the lower bound for the cost's Lipschitz constant. Our idea is thus to extend this algorithm to its META version, relying on the intuition that information collected solving similar problems can help in improving the choices of these initial parameters, ultimately enhancing the optimization procedure.

To this end, let us introduce the similarity vector

$$\mathcal{S} = \begin{bmatrix} \mathcal{S}_1 & \cdots & \mathcal{S}_M \end{bmatrix}^\top \in \mathbb{R}^M, \quad (19)$$

whose elements satisfy the following relationships:

$$\mathcal{S}_i \geq 0, \quad i = 1, \ldots, M, \quad (20a)$$

$$\sum_{i=1}^M \mathcal{S}_i = 1. \quad (20b)$$

This vector characterizes the relative similarity between the problem we aim at solving and the $M$ ones whose features are included in the META-dataset $\mathcal{D}^{\mathrm{meta}}$. Note that, if the new problem has already been solved (and it corresponds to the one associated with the $m$-th instance of the META-dataset), then we have $\mathcal{S}_m = 1$ and $\mathcal{S}_j = 0$, for all

$j \in \{1, \ldots, M\}$, $j \neq m$. Let us then define the META-initialization of $X^{(1)}$ and $\underline{\gamma}_f$ as follows:

$$X^{(1),\mathrm{meta}} = \sum_{i=1}^M \mathcal{S}_i X_i^{\star(n_{max})}, \quad (21a)$$

$$\underline{\gamma}_f^{\mathrm{meta}} = \sum_{i=1}^M \mathcal{S}_i \gamma_{f,i}^{(n_{max})}, \quad (21b)$$

so that the initial exploration point and lower bound on the Lipschitz constant for META-SMGO-$\Delta$ are constructed as convex combinations of the estimates of the global minima and Lipschitz constants comprised in $\mathcal{D}^{\mathrm{meta}}$.

Under the assumption that none of the $M$ instances of SMGO-$\Delta$ considered in the construction of $\mathcal{D}^{\mathrm{meta}}$ has been trapped in a local minima, we can now formalize the impact of the META-initialization on the difference between the initial estimate of the minimal function value and the true minimum as follows.

*Proposition 1:* Consider problem (1) and assume that its cost function $f(X)$ is $\rho$-similar and $\zeta$-similar (in the spirit of Definition 1) to a set of $M$ functions $\{f_i(X)\}_{i=1}^M$ for which (1) has already been solved without being trapped by a local minima. Assume that $\rho \leq v^{(1)}$, with $v^{(1)}$ characterizing the trust region $\mathcal{T}^{(1)}$ according to (13). Further assume that $X^{(1),\mathrm{meta}}$ in (21a) is a feasible initial exploration point. Under these assumptions, for the first exploration point obtained with a META-initialization (in (21a)-(21b)), the following bound holds

$$z_\theta^{(1)} - z^\star \leq 2\rho(\gamma_f^{\max} + \zeta), \quad (22)$$

where

$$z_\theta^{(1)} = \underline{f}^{(1)}(X_\theta^{(1)}), \quad \gamma_f^{\max} = \max_{i=1,\ldots,M} \gamma_{f,i}^{(n_{max})}. \quad (23)$$

*Proof:* Since the optimization subroutines used to populate the META-dataset are assumed to be $\rho$-similar, then the following holds

$$\left\| X_i^{\star(n_{max})} - X^\star \right\|_2 \leq 2\rho. \quad (24)$$

The distance between $X^{(1),\mathrm{meta}}$ and the optimal solution can thus be bounded as

$$\begin{aligned} \left\| X^{(1)} - X^\star \right\|_2 &= \left\| \sum_{i=1}^M \mathcal{S}_i X_i^{\star(n_{max})} - X^\star \right\|_2 \\ &\leq \sum_{i=1}^M \mathcal{S}_i \left\| X_i^{\star(n_{max})} - X^\star \right\|_2 \\ &\leq 2 \sum_{i=1}^M \mathcal{S}_i \rho = 2\rho \end{aligned} \quad (25)$$

where the second inequality holds thanks to the properties of the similarity vector $\mathcal{S}$ in (20) and the bound (24). As $\rho \leq v^{(1)}$ by assumption, then $X^* \in \mathcal{T}^{(1)}$ and thus, (14) holds. This implies that

$$z_\theta^{(1)} \leq z^{\star(1)} - \alpha \gamma_f^{(1)} \quad (26)$$

---

[2]This assumption is likely to be verified in many practically relevant applications where constraints are more general and less dependant on the specific context, e.g., the basic traffic rules in autonomous driving.

**Algorithm 1** Meta SMGO-$\Delta$

---

**Require:** $\mathcal{S}^{(1)}, \mathcal{D}^{\mathrm{meta}}$

$\quad X^{(1)} = \left(\mathcal{S}^{(1)}\right)^{\top} \left[ X_1^{\star(n_{max})} \quad \cdots \quad X_M^{\star(n_{max})} \right]^{\top}$

$\quad \underline{\gamma}_f^{(1)} = \left(\mathcal{S}^{(1)}\right)^{\top} \left[ \gamma_{f,1}^{\star(n_{max})} \quad \cdots \quad \gamma_{f,M}^{\star(n_{max})} \right]^{\top}$

$\quad$ **while** $n \leq n_{\max}$ **do**

$\qquad$ Evaluate $z^{(n)} = f(X^{(n)})$, $c_s^{(n)} = g_s(X^{(n)})$

$\qquad$ **for** $i = 1, \ldots, M$ **do**

$\qquad\quad$ Compute $\hat{z}_i^{(n)}$ as in (30)

$\qquad$ **end for**

$\qquad$ Find $\hat{\mathcal{S}}^{\star}$ by solving (31)

$\qquad$ Update $\mathcal{S}_{\theta}^{(n)} = \mathcal{S}^{(n)}$ as in (32)

$\qquad$ Find $X^{(n+1)}$ with modified SMGO-$\Delta$ (33)

$\quad$ **end while**

---

Subtracting on both sides the actual value $z^{\star}$ of the function we aim at optimizing at the global optimimum we further obtain:

$$
\begin{aligned}
z_{\theta}^{(1)} - z^{\star} &\leq z^{\star(1)} - \alpha\gamma_f^{(1)} - z^{\star} \\
&\leq \gamma_f \|X^{(1)} - X^{\star}\|_2 - \alpha\gamma_f^{(1)} \\
&\leq 2\rho\gamma_f - \alpha\gamma_f^{(1)} \leq 2\rho\gamma_f, \quad (27)
\end{aligned}
$$

where $\gamma_f$ is the actual (unknown) Lipschitz constant of the function we are optimizing and the third inequality stems from the definition of Lipschitz continuity. Adding and subtracting on the right-hand-side of the previous inequality $2\rho\gamma_f^{(1)}$ we further obtain:

$$
z_{\theta}^{(1)} - z^{\star} \leq 2\rho\gamma_f^{(1)} + 2\rho(\gamma_f^{(1)} - \gamma_f) \leq 2\rho(\gamma_f^{(1)} + \zeta). \quad (28)
$$

Since $\gamma_f^{(1)} = \underline{\gamma}_f^{\mathrm{meta}}$, based on the definition of $\underline{\gamma}_f^{\mathrm{meta}}$ in (21b), it straightforwardly follows that $\gamma_f^{(1)} \leq \gamma_f^{\max}$, thus concluding the proof. ∎

The previous bound holds for any $\mathcal{S}$ satisfying (20), yet it is of paramount importance for the similarity vector $\mathcal{S}$ to provide a reliable estimate of the similarity between the problem we aim at tracking and the $M$ ones that we have already solved to further reduce the distance of $z_{\theta}^{*(n)}$ from $z^*$. Toward this goal, in this paper we propose to iteratively evaluate similarity through the META-SMGO-$\Delta$ iterations as summarized in Algorithm 1, thus considering an iteration varying similarity vector.

At the beginning of the new optimization routine no information is available on the new function to be optimized. Therefore, we initially impose

$$
\mathcal{S}^{(1)} = \left[ \tfrac{1}{M} \quad \tfrac{1}{M} \quad \cdots \quad \tfrac{1}{M} \right]^{\top}, \quad (29)
$$

not to (wrongly) prioritize any instance of the META-dataset with respect to the others. Since at each new iteration $n \in [1, n_{max}]$ of META-SMGO-$\Delta$ we have access to a new evaluated point $X^{(n)}$, they are then incrementally employed to refine our initial guess. Specifically, we evaluate the unknown function we aim at optimizing at the current data point, namely

$$
z^{(n)} = f\left(X^{(n)}\right).
$$

To extrapolate the updated similarity vector $\mathcal{S}^{(n)}$, the latter is then compared with the following natural neighbor interpolation[3] of

$$
\hat{z}_i^{(n)} = \sum_{k=1}^{n_{max}} w_i(X_i^{(k)}) z_i^{(k)}, \quad i = 1, \ldots, M, \quad (30)
$$

where we use a set of Laplacian weights [17]. In particular, we solve the optimization problem (nested in the SMGO-$\Delta$ baseline routine):

$$
\min_{\hat{\mathcal{S}}} \quad \|z^{(n)} - \hat{\mathcal{S}}^{\top}\hat{Z}\|_2^2 \quad (31a)
$$

$$
\text{s.t.} \quad \hat{\mathcal{S}}_i \geq 0, \quad i = 1, \ldots, M, \quad (31b)
$$

$$
\sum_{i=1}^{M} \hat{\mathcal{S}}_i = 1, \quad (31c)
$$

where $\hat{Z} = \left[ \hat{z}_1^{(n)} \quad \cdots \quad \hat{z}_M^{(n)} \right]^{\top}$. Its solution $\hat{\mathcal{S}}^{\star}$ is then combined with the similarity vector available at the beginning of current iteration, namely

$$
\mathcal{S}_{\theta}^{(n)} = \mathcal{S}^{(n)} + \tau^{n-1}\hat{\mathcal{S}}^{\star}, \quad (32)
$$

where $\tau \in [0, 1]$ is a *discounting factor* introduced to promote smoothness in the similarity estimate over consecutive iterations. Clearly, (32) does not satisfy the properties of the similarity vector (see (20)). Accordingly, the elements of $\mathcal{S}_{\theta}^{(n)} = \mathcal{S}^{(n+1)}$ are then normalized, in order to lay in $[0, 1]$. This updated similarity matrix is used to refine $\underline{\gamma}_f^{(n)}$ as

$$
\underline{\gamma}_f^{(n)} = \left(\mathcal{S}_{\theta}^{(n)}\right)^{\top} \left[ \gamma_{f,1}^{(n_{max})} \quad \cdots \quad \gamma_{f,M}^{(n_{max})} \right]^{\top},
$$

*before* the regular SMGO-$\Delta$ Lipschitz constant estimation. Additionally, after the exploitation subroutine of SMGO-$\Delta$, the estimated similarity vector is also used to promote sampling near the updated estimate of the optimal value as follows:

$$
X_{\vartheta}^{(n)} = (1 - \tau^{n-1})X_{\theta}^{(n)} + \tau^{n-1}(\mathcal{S}_{\theta}^{(n)})^{\top}\mathbf{X}^{\star(n_{max})}. \quad (33)
$$

where $\mathbf{X}^{\star(n_{max})}$ is a vector stacking the estimated global minima comprised in $\mathcal{D}^{\mathrm{meta}}$. This new point is then tested with the *expected improvement* test (14), and selected as the new sampling point $X^{(n+1)}$ if it passes this check. Otherwise, the exploration routine takes place with no difference with respect the original SMGO-$\Delta$. Note that, in this case, the discounting factor is of fundamental importance to exploit meta information at the beginning of the optimization procedure, and gradually relying on the SMGO-$\Delta$ capabilities of finding the minimum when the number of iterations increases.

---

[3]Accordingly, the more iterations $n_{max}$ are performed, the more $\hat{z}_i^{(n)}$ will be informative.

| $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | $a_6$ | $a_7$ | $a_8$ |
|---|---|---|---|---|---|---|---|
| 1 | 16 | 5 | 1 | 16 | 5 | 2 | 80 |

TABLE II: Parameters for $\mathcal{D}^{\mathrm{meta}}$ generation

| n. iter | $\Delta$ | $\beta$ | $\alpha$ | $X^{(1)}$ | $\delta_{\max}$ | $\tilde{\rho}$ | $\tilde{\zeta}$ |
|---|---|---|---|---|---|---|---|
| 500 | 0.5 | 0.1 | 0.1 | (0.4775,0.0667) | 75% | 1.5 | 3000 |



Fig. 1: Contour plot of $f^o$. In grey the unfeasible areas due to $g_1$ and $g_2$. In green, the global (cross) and the local (circles) minima, while the blue ones the sampled $\{X\}_{n=1}^{n_{\max}}$.

## V. NUMERICAL EXPERIMENTS

In this preliminary work, the advantages of meta-learning methodology are illustrated with the low-dimensional ($X \in \mathbb{R}^2$) example taken from [16] (in noiseless settings). The objective functions belong to the class of the parameterized Styblinski-Tang function with offset, defined as:

$$f(X, \varphi_f) = \frac{1}{a_7}\Big(a_1 X_1^4 - a_2 X_1^2 + a_3 X_1 + \tag{34}$$

$$a_4 X_2^4 - a_5 X_2^2 + a_6 X_2 + a_8\Big),$$

with parameters $\varphi_f = [a_1, \cdots, a_8]$. Each $a_j, j = 1, \ldots, 8$ is randomly obtained as $a_j = a_j^o + \delta_j$, where $a_j^o$ is the nominal value as reported in Tab. I, and $\delta_j$ is a perturbation $\delta_j \sim \mathcal{U}(-1, 1) \cdot \delta_{\max}$. The constraints are fixed and equal to

$$g_1(X) = -4 + \big\| X - [-2.90, 2.90]^\top \big\|_2, \tag{35}$$

$$g_2(X) = \cos\big(2\big\| X - [\ 2.90, 2.90]^\top \big\|_2\big), \tag{36}$$

independently from the considered problem instance. The meta-data-set $\mathcal{D}^{\mathrm{meta}}$ (containing $M = 10$ meta-functions) is generated optimizing each $f_m \in \mathcal{D}^{\mathrm{meta}}$ with SMGO-$\Delta$, by using the parameters values reported in Table II. The generated functions are $\tilde{\rho} = 1.5$ $\tilde{\zeta} = 3000$ similar to each other. In Fig. 1, the contour plot of the nominal $f^o$ and the sampled points during the procedure are displayed, showing that the global minimum (green cross) is centered in $(-2.90, -2.90)$.

### A. Limit case: $f \equiv f_{\tilde{m}}$

The case where the new $f$ is equivalent to an already-optimized function $f_{\tilde{m}} \in \mathcal{D}^{\mathrm{meta}}$ is first tested, to verify
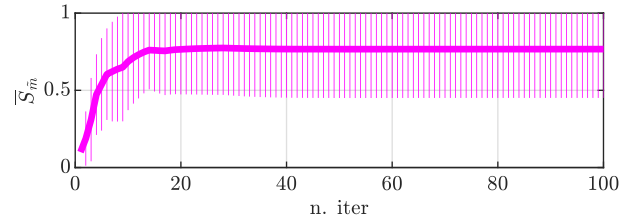


Fig. 2: Trajectory of the average coefficient values $S_{\tilde{m}}$, where $f = f_{\tilde{m}}$ (bold line) and their standard deviation (dictated by the envelope of the vertical lines).

that the data-driven algorithm identifies the correct similarity for $\mathcal{S}$. This holds true for a test with 10 repeated experiment, where the reference function changes such that $\tilde{m} = 1, \ldots, 10$. On average, the similarity coefficient $\mathcal{S}_{\tilde{m}}$ associated to the correct $f_{\tilde{m}}$ is 0.75, as shown in Fig. 2 and, consequently, $\mathcal{S}_i$ is small $\forall i \neq \tilde{m}$. This non-convergence to 1 is justified by the estimates of $\hat{Z}$ that are obtained with interpolation.

### B. General case

SMGO-$\Delta$ with meta-learning (META) is tested on $N = 10$ experiments[4], where each new objective $f_n, n = 1, \ldots, N$ is generated according to (34). Each $f_n$ is optimized with META ($\tau = 0.9$) and compared to standard SMGO-$\Delta$ optimization, with $n_{\max} = 100$ iterations. For fairness, in this preliminary work the Lipschitz constants of the constraints $\underline{\gamma}_{g_1}, \underline{\gamma}_{g_1} = 10^{-6}$ are initialized as if no meta-information is available, though it is clear that a prior that can be exploited to improve their initialization exists in the considered framework. Future work will be devoted to derive a more rigorous meta-formulation, that also considers this aspect linked to the constraints. Trajectories of the average *best* function value $\bar{z}^{*(n)}$ over the $N$ experiments obtained with META and SMGO-$\Delta$ are shown in Fig. 3, from which we can appreciate that META-SMGO-$\Delta$ significantly reduces the iterations required to reach the global minimum ($z_\theta^{(1)} - z^\star = 15.10 \leq 2 \cdot 1.5 \cdot (2817 + 3000)$, largely satisfying condition of Prop.1). In addition, the average number of infeasible samples is reduced by 25% (Fig. 5), even if no prior on the constraint is employed. Finally, notice how the initialization of $\underline{\gamma}_f$ is closer to the final estimate of $\gamma_f$ at convergence (Fig. 4), promoting a more targeted search of the regions where the minimum is expected from the start. This holds true also for functions where, due to the perturbation, the location of the (feasible) global minimum varies significantly from the nominal optimum (Fig. 5). In addition, in case the META-initialization of $\underline{\gamma}^{(1)}$ is an overestimate of the true $\gamma$, the update of $\mathcal{S}$ can compensate the error such that $\tilde{\gamma}^{(n)} \leq \tilde{\gamma}^{(n+1)}$.

### C. Sensitivity to M

The sensitivity to size $M$ of the meta-data-set is tested for $M = 5, 10, 20, 40$, on the $N = 10$ test experiments. On average, both the distance from the optimal value $z^{(1)} - $

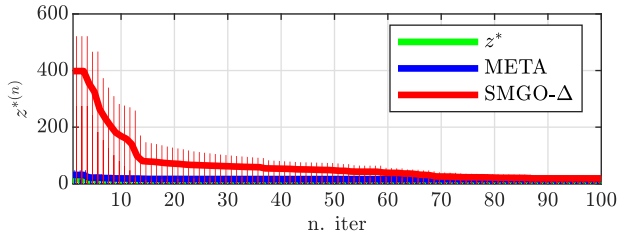[4]Accordingly, $\mathcal{S}^{(1)} = [0.1, \ldots, 0.1]^\top$

Fig. 3: Average trajectory of $z^{*(n)}$ during optimization and its standard deviation (envelope of the vertical lines).
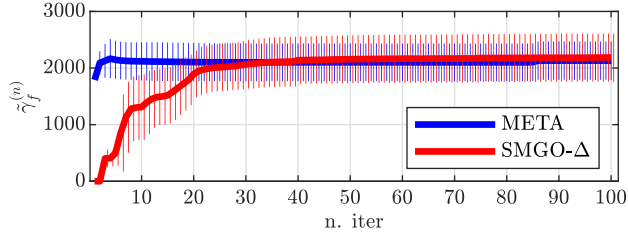


Fig. 4: Average trajectory of $\tilde{\gamma}_f^{(n)}$ during optimization and its standard deviation (envelope of the vertical lines).
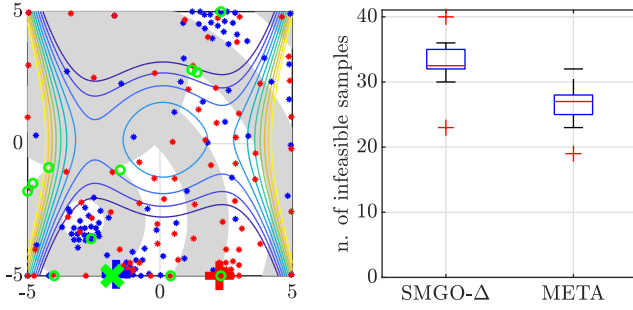


Fig. 5: Contour plot of $f$, with sampled points with META (blue) and SMGO-$\Delta$ (red) [left panel] and average number of infeasible points during the optimization [right panel].



Fig. 6: Sensitivity to $M$: average $z^{*(n)}$ [upper panel] and average number of constraint violations [lower panel] over $N = 10$ experiments.

$z^*$ and constraint violations (Fig. 6) are reduced for greater $M$. Nonetheless, for $M = 40$ there is a settling of these improvements. Numerical results thus confirm the intuition that a more examples are informative up to that point where they can eventually become redundant.

## VI. CONCLUSIONS

This work applies the meta-learning rationale to SMGO-$\Delta$, exploiting the similarity between optimization problems already solved with the one at hand to initialize two hyper-parameters of the nominal method, namely $X^{(1)}$ and $\underline{\gamma}_f$. We demonstrate that such an initialization results in a theoretical bound on the closeness of $z_\theta^{(1)}$ to the global minimum $z^*$. Numerical experiments confirm the theoretical findings and demonstrate that faster convergence and reduced constraint violations can be attained by relying on a meta-learning rationale. Further works will focus on the extension to non-fixed constraints and time-varying functions.

## REFERENCES

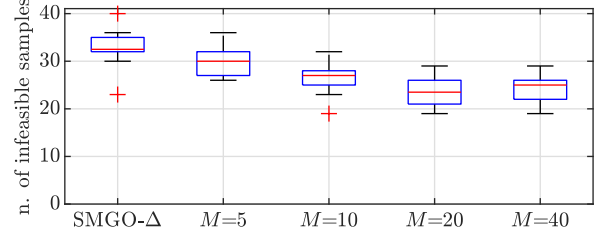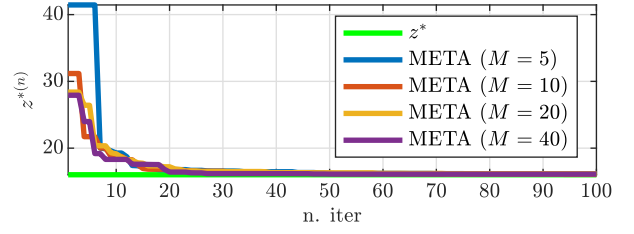[1] E. Brochu, V. M. Cora, and N. de Freitas, "A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning," *arXiv preprint arXiv:1012.2599*, 2010.

[2] J. Nocedal and S. J. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.

[3] S. Kirkpatrick, C. D. Gelatt Jr, and M. P. Vecchi, "Optimization by simulated annealing," *science*, vol. 220, no. 4598, pp. 671–680, 1983.

[4] J. Kennedy and R. C. Eberhart, "Particle swarm optimization," *Proceedings of the IEEE international conference on neural networks*, vol. 4, pp. 1942–1948, 1995.

[5] S. Thrun and L. Pratt, "Learning to learn," *Springer Science & Business Media*, 1998.

[6] R. Vilalta and Y. Drissi, "A perspective view and survey of meta-learning," *Artificial intelligence review*, vol. 18, no. 2, pp. 77–95, 2002.

[7] J. Van Schoren, B. Dorronsoro, and Y. Bengio, "Meta-learning: A survey," *arXiv preprint arXiv:1810.03548*, 2018.

[8] Y. Wang, W. Zhang, D. Ramanujan, and Y. Yang, "Generalizing meta-learning via task-aware modular architecture," *arXiv preprint arXiv:2007.04131*, 2020.

[9] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1126–1135, 2017.

[10] D. Li, Y. Yang, F. Ye, and T. M. Hospedales, "Episodic training for low-resource domain adaptation," *arXiv preprint arXiv:1902.09900*, 2019.

[11] S. M. Richards, N. Azizan, J.-J. Slotine, and M. Pavone, "Control-oriented meta-learning," *The International Journal of Robotics Research*, p. 02783649231165085, 2022.

[12] T. Guo, A. A. Al Makdah, V. Krishnan, and F. Pasqualetti, "Imitation and transfer learning for LQG control," *IEEE Control Systems Letters*, 2023.

[13] L. Ecker and M. Schöberl, "Data-driven control and transfer learning using neural canonical control structures," *arXiv preprint arXiv:2302.04042*, 2023.

[14] R. Busetto, V. Breschi, and S. Formentin, "Meta-learning for model-reference data-driven control," *arXiv preprint arXiv:2308.15458*, 2023.

[15] A. Chakrabarty, "Optimizing closed-loop performance with data from similar systems: A Bayesian meta-learning approach," in *2022 IEEE 61st Conference on Decision and Control (CDC)*. IEEE, 2022, pp. 130–136.

[16] L. Sabug Jr, F. Ruiz, and L. Fagiano, "SMGO-$\Delta$: Balancing caution and reward in global optimization with black-box constraints," *Information Sciences*, vol. 605, pp. 15–42, 2022.

[17] T. A. Bobach, "Natural neighbor interpolation-critical assessment and new contributions," Ph.D. dissertation, Technische Universität Kaiserslautern, 2009.