

Global Optimization for mixed categorical-continuous variables based on Gaussian process models with a randomized categorical space exploration step

M. Munoz Zuniga^a, D. Sinoquet^a

^aIFP Energies Nouvelles, 1-4 Avenue du bois préau 92852 Rueil-Malmaison Cedex - France

ARTICLE HISTORY

Compiled August 16, 2019

ABSTRACT

Real industrial studies often give rise to complex optimization problems involving mixed variables and time consuming simulators. To deal with these difficulties we propose the use of a Gaussian process regression surrogate with a suitable kernel able to capture simultaneously the output correlations with respect to continuous and categorical/discrete inputs without relaxation of the categorical variables. The surrogate is integrated into the Efficient Global Optimization method based on the maximization of the Expected Improvement criterion. This maximization is a Mixed Integer Non-Linear problem which is solved by means of an adequate optimizer: the Mesh Adaptive Direct Search, integrated into the NOMAD library. We introduce a random exploration of the categorical space with a data-based probability distribution and we illustrate the full strategy accuracy on a toy problem. Finally we compare our approach with other optimizers on a benchmark of functions.

KEYWORDS

Derivative Free Optimization; Surrogate Models; EGO; NOMAD; Categorical variables

1. Introduction

1 The field of research around mixed integer non-linear programming (MINLP) has
2 recently focused on designing algorithms specifically dedicated to finding global solu-
3 tions. Nevertheless, this latter task is all the more difficult than the model to optimize
4 is expensive to evaluate and thus often relies on the construction of a cheap-to-evaluate
5 surrogate. Hence, iterative surrogate-based approaches have been developed in the lit-
6 erature and can be decomposed, as presented in [Muller et al. \(2013\)](#), as

- 7 (1) Build an initial experimental design and evaluate the optimized function
- 8 (2) Compute the surrogate model based on the available evaluations
- 9 (3) Select the next sample point(s) with respect to some improvement surrogate-
10 based criteria
- 11 (4) Update the surrogate model with the new evaluate point(s)
- 12 (5) Iterate through (3) and (4) until a predefined stopping criterion has been met.

13 In this context, the Efficient Global Optimization (EGO) method ([Jones et al., 1998](#))

14 has given noticeable results, in particular, when dealing with expensive blackbox
15 simulators, as in Kanazaki et al. (2015), Hamza and Shalaby (2014), Comola et al.
16 (2016), for instance. EGO is based on a Gaussian process (GP) surrogate (Rasmussen,
17 2006) and an adaptive strategy where one or more new points are iteratively selected,
18 to be evaluated, with respect to the so called Expected Improvement (EI) criterion.
19 This criterion offers a trade-off between exploitation and exploration by adding
20 points around potential optima and unexplored areas. We propose to extend the
21 EGO strategy to mixed discrete-continuous inputs and will specifically focus on
22 the categorical case where no order is presumed on the discrete variables. To deal
23 with discrete variables we make use of a dedicated kernel proposed in the literature
24 of GP based surrogates (Zhou et al., 2011; Qian et al., 2008). We then integrate
25 this surrogate into the EGO strategy. Once the surrogate is able to deal with
26 categorical variables, the difficulty lies in the optimization of the EI criterion. This
27 latter sub-optimization problem is also a mixed continuous-discrete one, but with
28 an EI function that is relatively cheap to evaluate. Up to this point our strategy is
29 very similar with the recent work of Pelamatti et al. (2018). But instead of relaxing
30 the categorical variables into ordered integer or continuous variables, we use the
31 Mesh Adaptive Direct Search algorithm (Audet and J. E. Dennis, 2006) available
32 in the NOMAD library to solve the mixed categorical-continuous sub-problem.
33 This algorithm initially requires a notion of proximity in order to explore the input
34 space and can therefore straightforwardly deal with continuous and ordered integer
35 variables. In the presence of categorical variables the discrete space exploration is
36 left to the user in the NOMAD implementation. We take advantage of this latter
37 opportunity and develop a random exploration strategy of the discrete space, given
38 by a discrete probability distribution that can evolve as the optimization proceeds.

39
40 In comparison with the current literature that we are aware of, the contribution of
41 this article relies on three points. Firstly, we propose to deal, from the beginning
42 to the end, with both continuous and discrete variables in the surrogate based
43 optimization without relaxing categorical variables to integer or continuous ones.
44 Secondly we construct a discrete distribution used to explore the categorical vari-
45 ables randomly within the EI maximization with NOMAD. Lastly we propose a
46 comprehensive presentation and comparison with some well established radial basis
47 functions (RBF) surrogate-based approaches of the literature. Indeed most of the
48 literature on RBF-based optimization does not discuss the importance of kernels and
49 hyper-parameter choices and the key implication/simplification it can involve in the
50 problem formulation.

51
52 In section 2, we will introduce the main context and notations, followed in section 3
53 by a review of the GP surrogate approach: model, kernels and discussion. In section
54 4, the EGO algorithm framework adapted to the discrete case is presented, with a
55 focus on NOMAD algorithm for the Expected Improvement criterion optimization.
56 In section 5, we introduce a novel discrete probability distribution for the random
57 sampling of the categorical variables within NOMAD. In section ??, we present a
58 couple of surrogate-based methods for comparison and discuss the similarities and
59 differences with our approach. Finally in section 6 we apply the proposed methodology
60 to a benchmark of test functions and we compare our method with other surrogate
61 optimization methods based on RBF.

62

63 2. Context and notations

64 Our approach can be placed in the framework of "engineering models" optimization
65 as introduced in Swiler et al. (2014). Our purpose is to tackle optimization prob-
66 lems involving computationally expensive simulators with a moderate number of
67 optimization variables. Indeed, in the high-dimensional input cases, the number of
68 data points necessary to capture the function structure increases more rapidly when
69 categorical variables are involved; therefore, dealing with high-dimensional inputs
70 with expensive-to-evaluate simulators is rather difficult. We give a hint on the reason
71 why this difficulty arises in section 3.3 in the context of Gaussian process based
72 optimization.

73
74 The function to optimize will be denoted by f and the mixed parameters by $w =$
75 (x, z) where x represents the continuous variable vector of dimension p and z the
76 integer/categorical one with dimension q . The discrete vector z is supposed to be
77 defined on $I = \prod_{i=1}^q I_i$ where I_i is a finite discrete set. If z_i is categorical then no
78 order is pre-supposed and $I_i = \{1, \dots, m_i\}$ where m_i is the number of levels of the i -th
79 categorical variable. The integers 1 to m_i are simply representation of the levels. In
80 the integer case, I_i is defined as $I_i = \{a_1, \dots, a_{m_i}\}$ where the a_j 's are ordered integer
81 numbers such that $j \leq k$ implies $a_j \leq a_k$.

82 Our aim is to solve the following optimization problem:

$$\begin{aligned} & \underset{w \in \mathbb{R}^p \times I}{\text{minimize}} && f(w) \\ & \text{subject to} && x \in B \subset \mathbb{R}^p, \end{aligned} \tag{1}$$

83 where B defines bound constraints. In the sequel we suppose that an initial design of
84 experiment (DoE) is given: $w = \{w^1, \dots, w^{n_0}\}$, with $w^i \in \mathbb{R}^p \times I$ and the corresponding
85 responses $\mathbf{y} = (y_1, \dots, y_{n_0})$ such that $y_i = f(w^i)$. The DoE will be iteratively enriched
86 with respect to an optimization scope as explained in the next sections. The chosen
87 initial DoE, of size n_0 , is the concatenation of independent Latin hypercube samplings
88 (LHS) (McKay et al., 1979; Santner et al., 2003), with respect to the continuous
89 variables, each of size $n_{0,i}$. We will further specify the chosen $n_{0,i}$ in the numerical
90 section. This kind of DoE displays good properties for GP surrogates in a sample of
91 test cases presented in Swiler et al. (2014). The results are competitive with sliced-
92 LHS (Qian, 2012) which are theoretically an adequate option for mixed designs but
93 more difficult to obtain. For these reasons, we just settle for the independent LHS DoE.

94
95 We now describe the first step of the proposed methodology which consists of the
96 construction of a Gaussian process surrogate of f .

97 3. Gaussian process surrogate with mixed inputs

98 3.1. GP Surrogate Model

99 Gaussian process models are flexible and efficient surrogates of complex computer
100 codes. The popularity of GP stems among other things from the availability of the
101 prediction distribution estimate. For the purpose of optimization, the availability of
102 the prediction distribution estimation opens up the possibility to devise refinement
103 strategies based on some measure of improvement in the regions of interest. We will

104 discuss this latter point in the next section.

105 We now present a specific Bayesian approach of the GP regression, as presented in
 106 [Helbert et al. \(2009\)](#), that will be used in the numerical tests. The GP regression is
 107 based on the hypothesis that the function f is a realization of a Gaussian process $f_{\mathcal{G}}$
 108 defined by a linear regression trend $\mu : \mathbb{R}^p \times I \rightarrow \mathbb{R}$, a constant variance σ^2 and a
 109 correlation function $K_{\theta} : (\mathbb{R}^p \times I)^2 \rightarrow [-1, 1]$ with hyper-parameters θ . The constant
 110 variance and the hyper-parameters are assumed known at this stage. For the sake
 111 of simplicity, the trend regression term will only take into account the continuous
 112 variables such that

$$\mu(w) = \sum_{i=0}^l \beta_i h_i(x),$$

113 where, for $i = 0, 1, \dots, l$, h_i are known functions (that have been chosen by the user
 114 from his prior information on the function trend), β_i are random coefficients modeled
 115 with some improper prior distribution and l is a non-negative integer. A mixed variable
 116 trend could be considered but not treated in this work. The conditional random process

$$f_{\mathcal{G}}^c := f_{\mathcal{G}} \text{ knowing that } f_{\mathcal{G}}(w^1) = y_1, \dots, f_{\mathcal{G}}(w^n) = y_n$$

117 is then known to also be Gaussian. $f_{\mathcal{G}}^c$ has known mean μ_c and correlation function
 118 K_c such that

$$\mu_c(w) = h^T(x)\hat{\beta} + r^T(w)R^{-1}(y - H\hat{\beta}) \quad (2)$$

119 and, in particular, the prediction variance $\sigma_c^2(w) = \sigma^2 K_c(w, w)$ is given by

$$\sigma_c^2(w) = \sigma^2 \left[1 - r^T(w)R^{-1}r(w) + v^T(w)(H^T R^{-1}H)^{-1}v(w) \right] \quad (3)$$

120 where the correlation matrix R of the DoE is defined by $R_{ij} = K_{\theta}(w^i, w^j)$, $i, j =$
 121 $1, \dots, n$. The cross-correlation vector between the prediction and the observations is
 122 denoted by $r(w) = [K_{\theta}(w, w^i)]_{i=1}^n$, while H is the matrix defined by $H_{ij} = h_j(x^i)$,
 123 $1 \leq i \leq n$, $1 \leq j \leq d$ and $v(w) = H^T R^{-1}r(w) - h(x)$. The vector $\hat{\beta}$ is explicitly given
 124 by

$$\hat{\beta} = (H^T R^{-1}H)^{-1}H^T R^{-1}y. \quad (4)$$

Hence, for known hyper-parameters θ and variance σ^2 , the predictor is given by μ_c (2)
 as in [Sacks et al. \(1989\)](#).

We calibrate the hyper-parameters and the variance by maximizing the log-likelihood.
 This log-likelihood is the logarithm of the probability of observing the experimental
 data with our GP model parametrized by θ and σ^2 , *i.e.*

$$\mathcal{L}(\theta, \sigma^2) = \ln \left[\frac{1}{(2\pi\sigma^2)^{n/2}|R|^{1/2}} \exp \left(-\frac{1}{2\sigma^2}(y - H\hat{\beta})^T R^{-1}(y - H\hat{\beta}) \right) \right]$$

125 The first-order optimality conditions result in analytical formula for σ^2 as a function
 126 of θ , namely,

$$\sigma^2(\theta) = \frac{1}{n}(y - H\hat{\beta})^T R^{-1}(y - H\hat{\beta}). \quad (5)$$

127 This latter expression of $\sigma^2(\theta)$ is plugged in the log-likelihood. The "plugged-in" or
 128 "concentrated" log-likelihood then boils down to

$$\mathcal{L}(\theta, \sigma^2(\theta)) = -\frac{1}{2} [n \ln(\sigma^2(\theta)) + \ln(|R|) + n + n \ln(2\pi)] \quad (6)$$

129 and has to be maximized with respect to θ . The gradient of the "concentrated"
 130 log-likelihood is analytical so that this non-linear optimization problem is generally
 131 tackled with a multi-start BFGS algorithm (Roustant et al., 2012).

132

133 Finally, the surrogate of the objective function is the conditional mean μ_c given by
 134 (2) in which we have plugged-in the hyper-parameter solution of the log-likelihood
 135 (6) optimization. In the sequel, we designate this surrogate by \hat{f} . For the prediction
 136 variance we also use the version with the plugged-in optimal hyper-parameters in (3)
 137 denoted by $\hat{\sigma}_c$ in the following.

138

139 We will now give more insight on the importance of the correlation kernel choice
 140 and the nature of the θ hyper-parameters in the continuous-discrete mixed variables
 141 context.

142 3.2. Correlation kernel for mixed-inputs

143 In this section we do not intend to present in details the large amount of literature on
 144 the correlation kernel choice and its implications in GP. Our aim is to give a sufficient
 145 intuition of its importance and to present the kernel we selected.

146 We first notice in (2) that the prediction at any point w can be written as the sum of
 147 a trend term and a linear combination of $r(w) = K_\theta(w, w_i)$. Hence the GP predictor is
 148 deeply impacted by the kernel choice. The mixed kernel, defined in $(\mathbb{R}^p \times I)^2$, is typically
 149 constructed with the association of two separate kernels: one for the continuous part
 150 defined in $(\mathbb{R}^p)^2$ and another one for the categorical part defined in I^2 such that

$$151 \quad K_\theta(w, w') = K_{\theta_{\text{Cont}}}(x, x') \times K_{\theta_{\text{Cat}}}(z, z')$$

152 with $w = (x, z) \in \mathbb{R}^p \times I$ and $w' = (x', z') \in \mathbb{R}^p \times I$. The vectors θ_{Cont} and θ_{Cat} are the
 153 hyper-parameters associated with their respective kernels. These latter will be defined
 154 in next paragraphs.

155 **Correlation kernel for continuous variables.** The continuous kernel part is the
 156 standard product of 1-D correlation kernels such that

$$157 \quad K_{\theta_{\text{Cont}}}(x, x') = \prod_{i=1}^p K_{\theta_i}(x_i, x'_i)$$

158 with $\theta_{\text{Cont}} = (\theta_1, \dots, \theta_p)$. In our context these hyper-parameters are called correlation
 159 lengths.

160 For the continuous part, the degree of smoothness of the stationary GP surrogate is di-
 161 rectly linked to the degree of smoothness of the associated positive definite correlation
 162 kernel (Rasmussen (2006), section 4.1.1). Moreover, for a 1-D kernel, the correlation
 163 length is associated with a notion of regularity, which is defined in Adler (1981) as
 164 the mean number of up-crossings of a given level u by the GP (a continuous random
 165 process achieves an up-crossing of a given level when its values crosses the level from
 166 below). The smaller (higher) the correlation length the higher (smaller) is the mean
 167 number of up-crossings. In the numerical applications we selected the anisotropic sta-
 168 tionary Matern-5/2 correlation kernel which offers enough flexibility to adequately
 169 capture the variability of numerous objective function depending on the choice of the
 170 correlation lengths:

$$171 \quad K_{\theta_{\text{Cont}}}(x^i, x^j) = \prod_{k=1}^p \left(1 + \frac{\sqrt{5}|x_k^i - x_k^j|}{\theta_k} + \frac{5(x_k^i - x_k^j)^2}{3\theta_k^2} \right) \exp \left(-\frac{\sqrt{5}|x_k^i - x_k^j|}{\theta_k} \right).$$

172 Note that the correlation lengths, θ_{Cont} , do not depend on the categorical levels which
 173 by construction implies that the correlation lengths are similar for all levels. In order
 174 to limit the number of hyper-parameters (and therefore the potential number of black-
 175 box simulations required to assess them) we assume this independence between the
 176 continuous hyper-parameters and the categorical levels. A version with dependence
 177 between continuous and categorical parameters is presented in Qian et al. (2008) and
 178 used in Han et al. (2009) in a Bayesian context. Thus the continuous/categorical inde-
 179 pendence hypothesis could be relaxed but at a computational cost which we will avoid
 180 for the benchmark tests presented in this paper.

181 **Correlation kernel for categorical variables.** In order to treat the categorical
 182 variables, different types of correlation kernel can be constructed and recent works in
 183 the literature focus on the implications and the relevance of the choice of these kernels.
 184 Again, here our goal is only to give an understanding of the nature and importance
 185 for the model of the kernel for categorical variables. For more details we invite the
 186 reader to see for instance: Pinheiro and Bates (1996), Qian et al. (2008), Pinheiro and
 187 Bates (2009), Zhou et al. (2011), Zhang and Notz (2015), Roustant et al. (2018) and
 188 Pelamatti et al. (2018).

189 Hereafter, the i^{th} categorical level is the value taken by z_i , and a global-level, denoted
 190 by c , is defined as a set of values assigned to the vector of categorical variables z . The
 191 total number of global-levels is given by

$$N_{\text{GL}} = \prod_{i=1}^q m_i, \quad (7)$$

192 with m_i the number of levels of the i^{th} categorical variable. The most flexible approach
 193 is then achieved by choosing the correlation kernel, defined on the finite set I^2 , as the
 194 correlation matrix T whose N_{GL}^2 elements are the correlations between global-levels
 195 pairs. This correlation kernel is thus defined as

$$K_{\theta_{\text{Cat}}}(c_i, c_j) = T_{c_i, c_j}, \quad (8)$$

196 for two global-levels c_i and c_j . Since the hyper-parameters represent correlations be-
 197 tween global-levels, the matrix T must be unit diagonal, symmetric and positive def-
 198 inite. In total generality, the elements of the matrix T are the $N_{\text{GL}}(N_{\text{GL}} - 1)/2$ cat-
 199 egorical hyper-parameters. Another approach considers independently the correlation
 200 between levels for each categorical variable. This involves q correlation matrices $T^{(k)}$
 201 of size m_k^2 and the kernel is defined as

$$K_{\theta_{\text{Cat}}}(c_i, c_j) = \prod_{k=1}^q T_{c_i^k, c_j^k}^{(k)}, \quad (9)$$

202 where c_i and c_j are two global-levels with, respectively, k -component: c_i^k and c_j^k .
 203 In this case the number of hyper-parameters is reduced to $\sum_{i=1}^q m_i(m_i - 1)/2$. In
 204 both cases, the kernels (8) and (9) take values in $[-1, 1]$. The latter approach can
 205 be justified by an hypothesis on the underlying structure of the GP model; i.e., it
 206 is supposed to be a weighted sum of independent GPs with the same correlation
 207 function K_{Cont} (one GP per level), see Qian et al. (2008) for further details.

208
 209 The reduced number of hyper-parameters therefore comes with an underlying
 210 hypothesis on the structure of the GP which might not be adequate in some cases,
 211 in the sense that even the maximum log-likelihood solution θ could give a poor
 212 representation of the objective function if the underlying GP model structure is too
 213 far from the real function. In this case the approach (8) with N_{GL} hyper-parameters
 214 might be more appropriate if affordable. In the sequel we will use (9) as categorical
 215 kernel for computational cost reasons.

216
 217 To simplify the log-likelihood optimization task which is a difficult positive definite
 218 constrained optimization problem with respect to the correlation matrix, we adopt
 219 the spherical parametrization of the Cholesky decomposition of each matrix $T^{(k)}$ as in
 220 Zhou et al. (2011). This latter trick transforms the previous constrained log-likelihood
 221 optimization problem into a box constrained one that can be solved with a BFGS
 222 algorithm (Byrd et al., 1995). In the sequel, θ_{Cat} is the categorical hyper-parameter
 223 vector of size

$$N_{\theta_{\text{SC}}} = \sum_{i=1}^q \frac{m_i(m_i - 1)}{2} \quad (10)$$

224 composed of the spherical coordinates associated with the correlation matrices.

225 **3.3. Discussion on the dimension of the input variables and** 226 **hyper-parameters**

227 We could further try to reduce the number of correlation parameters to be estimated
 228 but this would imply further simplifications of the model. Since our aim is to use
 229 this model to approximate black-box functions, on which little information on
 230 regularity is known, we prefer to adopt a flexible kernel. In fact our objective is
 231 to find the best trade-off between the flexibility of the model (good approximation
 232 skills) and its estimation cost (number of experiments required to well determine the
 233 hyper-parameters). We will discuss these points in the following.

234

235 Kernels for categorical variables often boil down to the product of correlation factors
 236 between the different levels. To save computational expenses, one can assume the
 237 same correlation between all the levels or at least per group (group to be defined
 238 as proposed in [Roustant et al. \(2018\)](#)). In this paper, we decided to not introduce
 239 this kind of prior information. Nevertheless, the kernel flexibility comes with a
 240 price. Indeed, the selected categorical kernel (9) involves a number of correlations to
 241 estimate that can increase rapidly with the number of levels as shown in (10). The
 242 more hyper-parameters we introduce, the greater is the flexibility of the surrogate
 243 model but in return more data points are required to capture enough information to
 244 "feed the flexibility". In this context, we will limit ourselves to applications with a
 245 few categorical variables with moderate number of levels. A large range of industrial
 246 problems falls in this framework ([Swiler et al., 2014](#)). Indeed, prior knowledge of
 247 mechanical engineers is often used to limit the number of possible levels of the
 248 categorical variables. For instance, in optimal design of a mechanical system, a few
 249 types of materials, predefined shapes or structures are selected beforehand.

250

251 4. Efficient Global Optimization with mixed inputs

252 4.1. Global optimization based on the Expected Improvement criterion

253 Once the hyper-parameters are tuned on the current DoE, the GP model is completely
 254 defined and can be used in an adaptive optimization scheme. The Efficient Global
 255 optimization strategy ([Jones et al., 1998](#)) relies on the posterior distribution of the GP
 256 model which enables us to assess the distribution of the following random improvement
 257 of f minimization:

258

$$I(w) = f_{min} - f_G^c(w) \quad (11)$$

259 and the related Expected Improvement (EI)

$$EI(w) = \mathbb{E}(\max(0, I(w))) \quad (12)$$

260 where $f_{min} = \min(y_1, \dots, y_{n_k})$ and n_k the size of the DoE at the k -th iteration of the
 261 method. This criterion gives a measure of the expected improvement achievable at a
 262 new point w (i.e. expectation to go below the current minimum) based on the known
 263 responses at the known available simulated points w_1, \dots, w_{n_k} and the GP surrogate
 264 estimated distribution. The Gaussian hypothesis enables the implementation of the
 265 following closed formula of the EI criterion ([Schonlau, 1997](#))

$$EI(w) = (f_{min} - \hat{f}(w))\Phi\left(\frac{f_{min} - \hat{f}(w)}{\hat{\sigma}_c(w)}\right) + \hat{\sigma}_c(w)\phi\left(\frac{f_{min} - \hat{f}(w)}{\hat{\sigma}_c(w)}\right) \quad (13)$$

266 where ϕ is the standard univariate Gaussian distribution and Φ its cumulative
 267 distribution function. This criterion offers a built-in exploration-exploitation measure
 268 for the optimization strategy.

269

270 We can now define the next point to simulate in our optimization scheme as the
271 point that maximizes the EI criterion. After evaluating the new point we can update
272 the GP model with the same hyper-parameters or update the θ 's too by maximizing
273 the updated log-likelihood. In the numerical results we will update the GP and the
274 hyper-parameters at each step.

275

276 At this stage, we can mention the work of [Taddy et al. \(2009\)](#) and [Gramacy and Taddy](#)
277 [\(2010\)](#) where a Tree GP (TGP) strategy is presented. On a tree structure, they propose
278 to construct a GP surrogate for each global-level independently: corresponding to the
279 leaves of the tree. The EI criterion (only depending on the continuous variables) is
280 evaluated on a sampled grid and the locations ranked with respect to their EI values
281 in an iterative manner. This seems to us rather costly and does not take into account
282 any correlation between the levels. Indeed a dense grid, per global-level, has to be
283 evaluated at each step. Nevertheless, this strategy is interesting as it enables to select,
284 from the ranking, not only one point but a batch of points to be evaluated at each
285 iteration. Adding only one point at the time would be equivalent to a simple grid
286 search per global-levels. We believe that doing one optimization per global-level would
287 give better results. This latter "per level" optimization strategy will be evaluated in
288 the numerical section and as mentioned will be considered equivalent or better than a
289 TGP with one point added at the time. Since our actual implementation of EGO does
290 not add batches of points, we will not compare it to the corresponding batch-TGP in
291 this paper.

292 *4.2. NOMAD for the EI criterion optimization*

293 The maximization of the EI criterion is also a mixed continuous-discrete problem
294 but with an objective function relatively cheap to evaluate. Sampling strategies are
295 sometimes preferred to optimize the improvement criterion ([Muller et al., 2013](#)) but
296 this seems inefficient since the parameter space to explore can be very large especially
297 with categorical variables.

298 To tackle this EI-maximization task we used the derivative free Mesh Adaptive Direct
299 Search (MADS) algorithm ([Audet and J. E. Dennis, 2006](#)). MADS is a robust opti-
300 mization method which can be used on a very wide range of non-linear optimization
301 problems. Nevertheless, when the direct search algorithm is not coupled with a surro-
302 gate model, the number of simulations required to reach the optimum can sometimes
303 be impractical when dealing with expensive-to-evaluate black-box models. For these
304 reasons we decided not to use MADS on the main optimization problem (1) but it
305 appears as a good tool for the EI maximization sub-problem.

306 We will now briefly describe the MADS algorithm and introduce our randomized
307 approach for the categorical space exploration in the next section.

308 The MADS algorithm consists of iteratively evaluating new trial points on an
309 adaptive mesh. Each iteration is divided into two steps: the search and the poll
310 steps. In the search step, a given number of trial mesh points are evaluated around
311 the current best point. If a better point (smaller value of the function) is found
312 the mesh is coarsened, if no improvement is achieved, the poll step is invoked. In
313 the poll step, new points to evaluate are chosen along random positive spanning
314 directions within a limited distance of the current best point. This distance is
315 controlled by a poll size parameter which is greater than or equal to the current
316 mesh size. This latter choice enables a possible dense exploration within the area

317 centered on the current best point and ensures convergence of the algorithm under
 318 mild hypotheses on the objective function (Audet and J. E. Dennis, 2006). Then,
 319 (Abramson et al., 2009) added to the MADS algorithm the notion of discrete
 320 neighbourhood, introduced in (Audet and Dennis, 2000). Also an additional extended
 321 poll step is introduced and triggered when no improvement is found in the two
 322 previous steps: a poll step is then performed around each point associated with
 323 an objective function value close enough to the current best one. This additional
 324 step can, in practice, help the algorithm to escape from some local optima. The
 325 described method is implemented in the NOMAD software and offers the option of a
 326 user-defined neighbourhood notion for the categorical variables for the poll step. We
 327 will use this capacity to define a probability based notion of proximity in our approach.
 328

329 **5. Random sampling of categorical variables within NOMAD for EI** 330 **optimization**

331 Often in the literature, categorical variables are coded as integers and then treated
 332 as real numbers or ordered integers. In this way the notion of neighbourhood is
 333 straightforward: often based on some l^p norm. In other cases categorical variables are
 334 coded as binaries (Potdar et al., 2017). In the binary space a notion of proximity is
 335 based on the number of flips necessary to pass from one global-level to another one
 336 (Hamming distance). These approaches can be useful when the user is able to define
 337 an imposed order on the categories or if the "binary flipping" proximity model has a
 338 real physical meaning. If this kind of information is not available, the ordering or/and
 339 the notion of proximity is clearly arbitrary and might skew the exploration in the
 340 optimization.

341
 342 In this context it seems natural to assume that, without any prior information avail-
 343 able, no proximity assumption should be introduced. A better approach seems to
 344 model the categorical variables as random variables with a discrete probability on the
 345 global-levels. So that, at each stage k of the optimization process of (1) (n_k points
 346 have been evaluated), the probability for the randomized categorical vector Z_k to take
 347 the global-level c_i is $p_{k,i}$ for $i = 1, \dots, N_{GL}$.

348 Within NOMAD, the aim of the extended poll step is to propose a new categorical
 349 position from the current one. Since no proximity assumption between categorical
 350 variables is considered, we suppose that the proposed random position, Z_k^{Poll} , should
 351 be chosen independently from the current one: $Z_k^{Current}$. But Z_k^{Poll} has to be different
 352 from the current position. So that the probability of switching from one current global-
 353 level to another in the poll step is given by

$$\mathbb{P}(Z_k^{Poll} = c_i | Z_k^{Current} = c_j, Z_k^{Poll} \neq c_j) = \frac{p_{k,i}}{1 - p_{k,j}} \quad (14)$$

354 for $i \neq j$ and 0 otherwise. Hence, at each step in the NOMAD algorithm, the user-
 355 defined neighbor of the current categorical variable is determined by sampling a global-
 356 level poll point according to the probability (14). Multiple global-level-polls points
 357 could be sampled similarly.

358 A first and natural idea to model the discrete probabilities $(p_{k,i})_{i=1, \dots, N_{GL}}$ is to consider

359 a uniform distribution so that

$$\mathbb{P}(Z_k^{Poll} = c_i | Z_k^{Current} = c_j, Z_k^{Poll} \neq c_j) = \frac{1}{N_{GL} - 1}. \quad (15)$$

360 In this case all global-levels (different from the current one) have the same probability
 361 to be chosen in the extended poll step of iteration k . In order to take advantage of
 362 the information on the distribution of the objective function evaluations on each
 363 global-level, we propose another approximation of $p_{k,i}$ in (14) based on a non-uniform
 364 discrete probability on the global-levels. We know, in particular, which global-levels
 365 have been more or less explored and which global-level is associated with the smallest
 366 objective function values. This information is already integrated within the EI
 367 function definition and its maximization should provide us the best next points in
 368 unexplored areas or areas where minimal objective function values are expected.
 369 Nevertheless, the extended poll step based on (15) does not explicitly use this
 370 information. We thus propose to integrate, in the extended poll step, the information
 371 learned from simulations of previous optimization iterations. This leads us to explore
 372 the global-levels randomly with respect to an "informative" probability distribution
 373 defined hereafter.

374

375 To give a hint on the relevance of the proposed discrete distribution, we can study the
 376 improvement criterion

$$377 \quad I_{Cat}(z) = \max(0, f_{\min} - M(z)),$$

378 with $f_{\min} = \min(y_1, \dots, y_{n_k})$ and

$$379 \quad M(z) = \min_{x \in B} f_G^c(x, z)$$

380 and $z \in \{c_1, \dots, c_{N_{GL}}\}$. For each z , $M(z)$ is a real random variable with cumulative
 381 distribution Ψ_z ¹, mean $\bar{M}(z)$ and standard deviation $\sigma_M(z)$. The expectation of this
 382 criterion can be developed as

$$\mathbb{E}(I_{Cat}(z)) = (f_{\min} - \bar{M}(z))\Psi_z(f_{\min}) + \sigma_M(z)\mathbb{E}\left[\frac{M(z) - \bar{M}(z)}{\sigma_M(z)}\mathbb{1}_{M(z) \leq f_{\min}}\right], \quad (16)$$

383 where $\mathbb{1}$ stands for the indicator function. Unfortunately, the quantity (16) can not
 384 be further computed since the distribution of $M(z)$ is unknown and too costly to ap-
 385 proximate empirically. Nevertheless, we will see that the proposed discrete probability
 386 presents some similarities with the terms involved in (16). We also should keep in mind
 387 that the prior information we hope to integrate in the EI exploration (13) has to be
 388 available at a small computational cost at least smaller than an intensive sampling
 389 of EI function, which can become cumbersome when the number of global-levels in-
 390 creases.

391 We propose that the discrete distribution on the global-levels integrates the density
 392 of evaluated points at each global-level. To achieve this goal, we introduce $p_{k,i}^g$ as the
 393 probability that the global-level c_i has not been fully explored. The quantity $p_{k,i}^g$ should
 394 thus be close to one when the global-level c_i needs more exploration. This quantity is

¹ Ψ_z is not necessary Gaussian

395 similar to the term $\sigma_M(c_i)$ in (16) which is large when uncertainty for the global-level
 396 c_i is large. Furthermore, the discrete distribution should integrate the potential of each
 397 global-level to contain the global minimum. For this purpose we introduce $p_{k,i}^m$ as the
 398 probability that the global-level c_i has a high potential of containing the minimum.
 399 Therefore, $p_{k,i}^m$ will present some similarity with the term $\Psi_z(f_{min})$ involved in (16).
 400 Finally, $p_{k,i}$, the probability of global-level c_i , is given by a weighted sum of the two
 401 previously introduced probability distributions, that is

$$p_{k,i} = \mathbb{P}(Z_k = c_i) = \alpha_k p_{k,i}^g + (1 - \alpha_k) p_{k,i}^m \quad \forall i = 1, \dots, N_{GL} \quad (17)$$

402 where k stands for the optimization iteration and $\alpha_k \in [0, 1]$ is a weight parameter.
 403 As k increases, the discrete distribution $p_{k,i}^g$ should converge to the zero discrete
 404 distribution since all global-levels will be fully explored. The probabilities $p_{k,i}^m$ should
 405 converge to a discrete distribution that is zero for all i except for the one associated
 406 with the global-level containing the global minimum.

407
 408 We define hereafter approximations of $p_{k,i}^g$ and $p_{k,i}^m$. The probability that the global-
 409 level c_i is not fully explored is approximated by

$$\hat{p}_{k,i}^g = 1 - \left(\frac{n_{k,i}}{n_k} \right)^l, \quad (18)$$

410 with $n_{k,i}$ the number of evaluated points in the global-level c_i at iteration k , n_k the
 411 total number of evaluations of the objective function at iteration k and $l > 0$ ($l = 1/2$
 412 in the numerical results). Then, the probability $\hat{p}_{k,i}^g$ depends on the proportion of
 413 points currently evaluated in the corresponding level. Selecting a decreasing sequence
 414 α_k in (17), with respect to k , will force $\alpha_k \hat{p}_{k,i}^g$ to go towards zero, and in a sense to
 415 mimic $\sigma_M(c_i)$ in (16).

416 The probability that the global-level c_i contains the minimum is approximated by

$$\hat{p}_{k,i}^m = \frac{S_{k,i}^R}{\sum_{j=1}^{N_{GL}} S_{k,j}^R}, \quad (19)$$

417 with $S_{k,i}^R$ an approximation of $\Psi_z(f_{min})$ in (16) for $z = c_i$. The chosen model for $S_{k,i}^R$
 418 is detailed in appendix A. It is defined as a function of the mean of the function values
 419 available for global-level c_i and the associated standard deviation.

420 Hence, NOMAD extended poll provides a global-level sampled according to (14)
 421 with $p_{k,i}$ given by (17) and $p_{k,i}^g, p_{k,i}^m$ respectively by (18) and (19).

422
 423 To summarize, the randomized approach for the extended poll step is based on
 424 available evaluated objective function data (initial DoE and from previous iterations),
 425 and is a trade-off between focusing on under-explored global-levels and the ones
 426 with potential optimality. The selected discrete distribution can be seen as a prior
 427 information integrated in the extended poll step. In the uniform case (15), the prior is
 428 non-informative, and we expect the NOMAD optimization to converge asymptotically,
 429 since all levels will be fully explored with probability one. When the proposed poll
 430 step proposition is given by (14) combined with (17), as described, we expect an
 431 accelerated recovery of the optimal EI which is illustrated by our numerical tests. The
 432 proof of convergence of the NOMAD algorithm with this latter extended poll scheme

433 seems tricky and dependent on the tuning parameters α_k , $b_{k,i}$ and $\hat{\sigma}_{k,i}$. Nevertheless,
 434 we are confident that it will at least converge to a local minimum which seems
 435 sufficient for the iterative (with respect to k) EI optimization. Convergence analysis
 436 will be the subject of a further work.

437

438 6. Numerical results

439 6.1. Two competitive approaches

440 In the numerical part we will compare our approach, denoted by Cat-EGO, to two
 441 different methods based on a radial basis function (RBF) surrogate: RBFOpt and
 442 MISO-CSTV(f). We will refer to the latter one as MISO. For details on the two algo-
 443 rithms we refer the reader to the corresponding papers, [Gutmann \(2001\)](#); [Costa and](#)
 444 [Nannicini \(2018\)](#) for RBFOpt and [Muller \(2016\)](#) for MISO.

445 Hereafter, we list the main differences between these two methods and our approach.
 446 First of all, MISO performs a local continuous optimization with categorical variables
 447 fixed at the values corresponding to the current best point. Concerning RBFOpt, an
 448 automatic selection of the kernel within the predefined set \mathcal{K} is performed by using
 449 a cross-validation scheme. One last significant difference between the MISO/RBFOpt
 450 strategies and ours is the initial DoE size. It is of the order of magnitude of $2(p+q+1)$
 451 for MISO and RBFOpt. This value corresponds to an approximation of the minimal
 452 number of data points required for the surrogate to be well fitted. Since our kernel
 453 has much more hyper-parameters, an adequate size of the initial DoE has to be much
 454 larger in order to sufficiently feed the surrogate model learning stage. Our DoE size
 455 is generally between $3 \times p \times N_{\theta_{SC}}$, with $N_{\theta_{SC}}$ given by (10), and $3 \times p \times N_{GL}$ with
 456 N_{GL} given by (7). This latter size gives more robustness to the hyper-parameter op-
 457 timization but becomes rapidly prohibitive if the problem involves a large number of
 458 global-levels N_{GL} .

459 The differences in DoE size appear as a consequence of the treatment of the categor-
 460 ical variables as continuous by both MISO and RBFOpt. This offers a much simpler
 461 function basis approach but also drastically reduces the potential to learn information
 462 within categorical global-levels. More precisely, no correlation is estimated between
 463 the global-levels. This correlation information helps for the global-level exploration:
 464 indeed, if two global-levels are detected as strongly correlated, then only one has to
 465 be explored. Another consequence of treating categorical variables as continuous is
 466 that a continuous interpolation between the global-levels is done: it assumes that the
 467 underlying regularity of this hypothetical continuous approximation of the categor-
 468 ical variables can be captured by the selected radial basis function. Moreover, the
 469 radial basis function model implicitly assumes stationarity of the approximated func-
 470 tion with respect to all variables, which is a very strong assumption when imposed on
 471 the categorical variables (when treated as continuous).

472 **6.2. A baseline function: exploration improvement for categorical**
473 **variables**

We consider a two-dimensional toy problem with one categorical variable with 10 levels defined as

$$f(x, z) = \begin{cases} \cos(3.6\pi(x - 2)) + x - 1 & \text{if } z = 1, \\ 2 \cos(1.1\pi \exp(x)) - \frac{x}{2} + 2 & \text{if } z = 2, \\ \cos(2\pi x) + \frac{1}{2}x & \text{if } z = 3, \\ x(\cos(3.4\pi(x - 1)) - \frac{x-1}{2}) & \text{if } z = 4, \\ -\frac{x^2}{2} & \text{if } z = 5, \\ 2 \cos(\frac{\pi}{4} \exp(-x^4))^2 - \frac{x}{2} + 1 & \text{if } z = 6, \\ x \cos(3.4\pi x) - \frac{x}{2} + 1 & \text{if } z = 7, \\ x(-\cos(7\frac{\pi}{2}x) - \frac{x}{2}) + 2 & \text{if } z = 8, \\ -\frac{x^5}{2} + 1 & \text{if } z = 9, \\ -\cos(5\frac{\pi}{2}x)^2 \sqrt{x} - \frac{\ln(x+0.5)}{2} - 1.3 & \text{if } z = 10. \end{cases}$$

474 This problem has several local minima with close function values (see Figure 1). Some
475 correlations between the individual one-dimensional functions associated with given
476 levels can be observed: *e.g.* functions at levels 7 and 10 are strongly correlated whereas
477 functions at levels 4 and 7 are anti-correlated.

478 NOMAD implementation for integer variables is compared on the minimization of
479 this function with the two proposed randomized poll steps for categorical variables:
480 uniform sampling and improved sampling which takes into account the current
481 simulated data distribution. The methods are run 100 times with different initial
482 points (the best of 5 randomly sampled points) in order to measure the robustness
483 of the methods. Figure 2 and 3 illustrate the effect of randomized categorical poll
484 steps on NOMAD results: the percentiles of runs reaching the global minimum area
485 increases from 40% to 60% (Figure 3) and the associated objective functions are
486 closer (see outliers in Figure 2 (right)). In comparison to the uniform strategy, the
487 improved approach shows more robust results with an higher rate of global minima
488 discovery.

489
490 In a second analysis, we compare the three NOMAD implementations on a "real"
491 EI function optimization. For this aim, we ran 800 times the EGO algorithm up to
492 the tenth iteration: EGO is run 80 times from random design of 5 experiments and
493 each iteration is performed 10 times to take into account the randomness of the EI
494 minimization methods. We present the results obtained for the next optimization (for
495 each of the 800 runs) of the EI function (the eleventh iteration) with the three NOMAD
496 implementations and two "per level" optimizations. The last two methods consist of
497 continuous optimizations with NOMAD and a multi-start BFGS method with fixed
498 categorical variable. For the mentioned example, ten 1D continuous optimizations are
499 run and the best EI of the 10 runs is considered as the solution. For NOMAD methods
500 applied on the mixed continuous-categorical space, 30 individual initial points are
501 uniformly sampled in the full space and NOMAD is run from the best point (minimal
502 EI). NOMAD-per-level is initialized by the best point of 3 uniformly sampled points
503 for each level and BFGS is run 3 times from 3 uniformly sampled initial points per
504 level.

505 Figure 4 displays the absolute EI errors of the solutions of the 5 methods compared
506 to the best EI among the solutions of all methods. We observe that the EI errors

507 obtained by the 5 methods are very close. NOMAD-per-level is more robust with a
508 very small errors for all runs, whereas the BFGS per level approach is less robust.
509 The latter method might converge to local solutions, whereas, NOMAD is a more
510 global method. On the other hand, in Figure 5 we observe that the number of
511 iterations necessary to reach the global EI is much larger for NOMAD-per-level
512 and NOMAD-integer, and slightly larger for the BFGS method compared to the
513 2 NOMAD methods with the randomized poll steps for categorical variables. The
514 accuracy of the strategies that consider each level independently comes with an
515 higher number of simulation cost. On this example, no significant differences are
516 noticeable between the uniform and the improved strategies. The same analysis has
517 been achieved with different EI-shaped functions associated with different iterations
518 of EGO, and similar conclusions are obtained on these cases.

519
520 Since the EI maximization is only a sub-optimization problem of the global cat-EGO
521 strategy, we now present the results of Cat-EGO on the toy example for the EI
522 maximization with the two proposed randomized NOMAD implementations and the
523 two "per-level" approaches with NOMAD and multi-start BFGS. The Cat-EGO
524 results are compared for 2 simulation budgets of 40 and 50. First, in Figures 6 and
525 7, we observe that the random sampling based strategies (Uniform and Improved
526 sampling) are more robust with regard to reaching the global optima with high
527 accuracy (0.001) for all the stopping criteria. The lack of robustness of multi-start
528 BFGS (for high accuracy) can be explained by the local optimization approach
529 of BFGS on the EI function which makes difficult "exploitation" of Cat-EGO. To
530 improve the exploration in this method, we should add more initial points for the
531 multi-start approach but it will become cumbersome for higher dimensions. On the
532 other hand, after a sufficient number of iterations (as for 40 and 50 simulations), the
533 multi-start BFGS strategy is very efficient in finding the settings of the categorical
534 variables at optimality (see percentages of success with accuracy 0.1: 85% and 90%).
535 NOMAD per level is less efficient in reaching the "right" level (74% and 78%) but is
536 more accurate in the optimization in continuous variable due certainly to its global
537 optimization skills (68% and 75% for accuracy 0.001 compared to 49% and 63% for
538 multi-start BFGS approach). In comparison to the two random sampling NOMAD
539 strategies, the 2 "per-level" EI maximization approaches are less efficient in reaching
540 the global optimum with high accuracy for the 2 simulation budgets: from 72% to
541 86% for accuracy 0.001.

542 The uniform sampling approach gives better results than the improved sampling
543 approach when the number of simulations is large enough: for a budget of 50
544 simulations, its percentage of success for the two accuracies are larger. The improved
545 sampling results in better percentages of success during the first iterations (see results
546 for the budget of 40 simulations) but seems to bias approximatively 5% to 7% of the
547 100 runs, leading to a smaller percentage of success in comparison with the uniform
548 sampling after 50 simulations. We remind that the tuning parameter α_k (17) was set
549 to zero in the numerical results and we expect a better behavior from the improved
550 strategy with an adequate adaptive tuning of α_k , increasing the exploration of levels
551 of the categorical variables compared to exploitation when necessary (the uniform
552 sampling performs only exploration).

553
554 In the next section the Cat-EGO will always be run with the NOMAD improved
555 random sampling poll step.

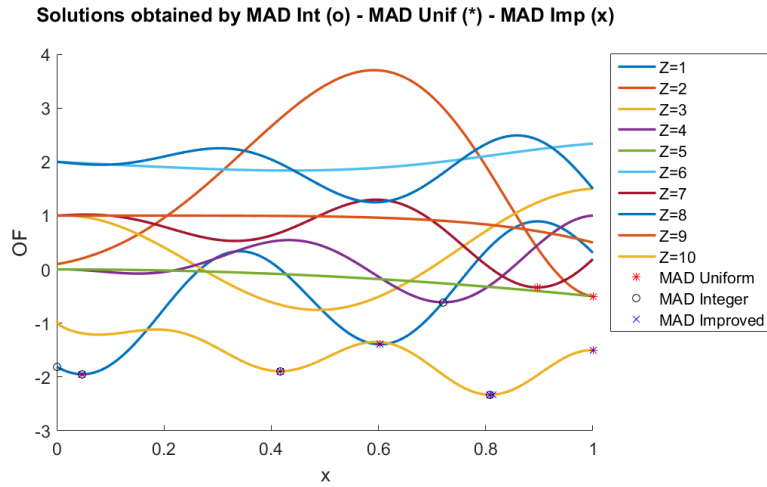


Figure 1. Toy problem with one continuous variable and one categorical variable with 10 levels. The global minimum is located at $x = 0.808$ on level 10 of the categorical variable, the associated objective function value is -2.329 . Solutions obtained by 100 runs of NOMAD with 3 different poll step strategies (adapted to integer, uniform sampling, "improved" sampling).

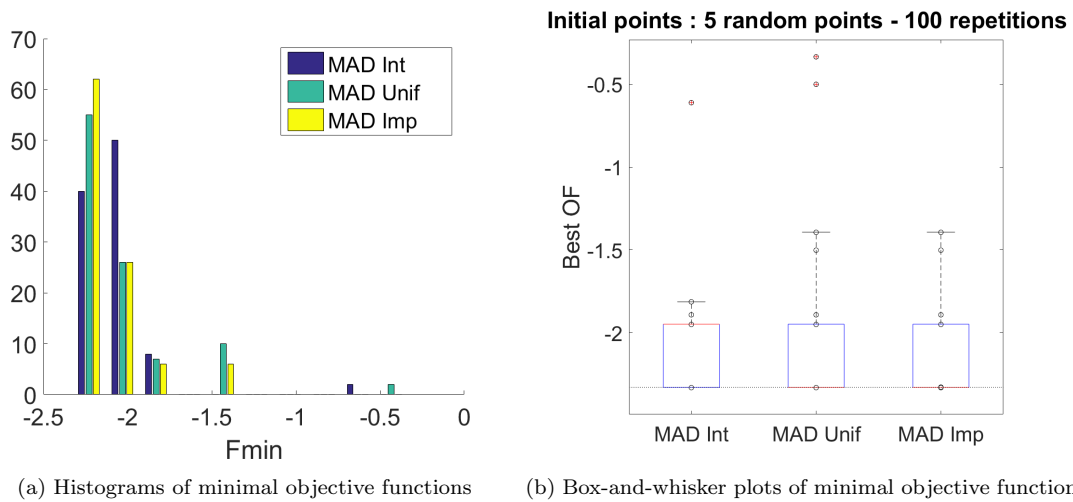


Figure 2. Minimal objective functions for 100 runs of NOMAD with 3 different poll step strategies (adapted to integer, uniform sampling, "improved" sampling). (a) Histograms of minimal objective functions. (b) Box-and-whisker plots of minimal objective functions, red lines indicate the medians (middle quartiles), the boxes include 50% of the values, the whiskers cover 99.3% of the values (under Gaussian distribution assumption), the red crosses being considered as outliers. The circles are the 100 minimal objective functions.

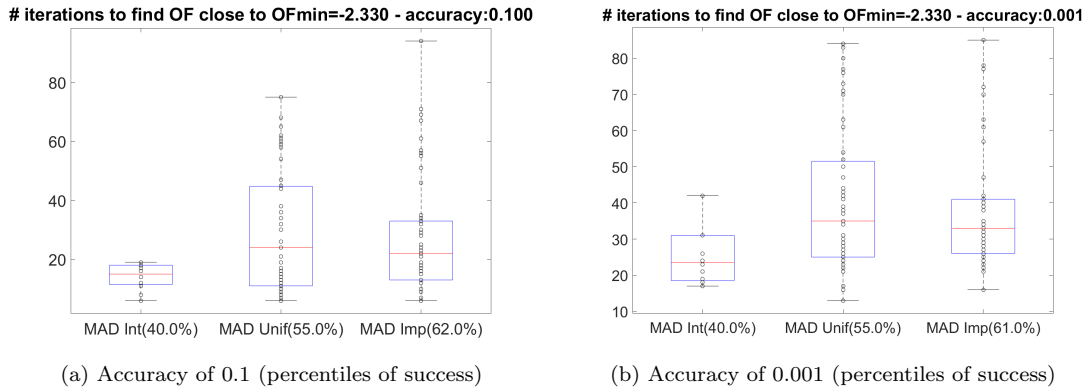


Figure 3. Box-and-whisker plots of the number of simulations necessary to reach the minimal objective functions with two given accuracies (0.1 and 0.001) for 100 runs of NOMAD with 3 different poll step strategies (adapted to integer, uniform sampling, "improved" sampling). The accuracy is on the absolute error of minimal objective function compared to global optimum. See legend of Figure 2 for details on Box-and-whisker plots.

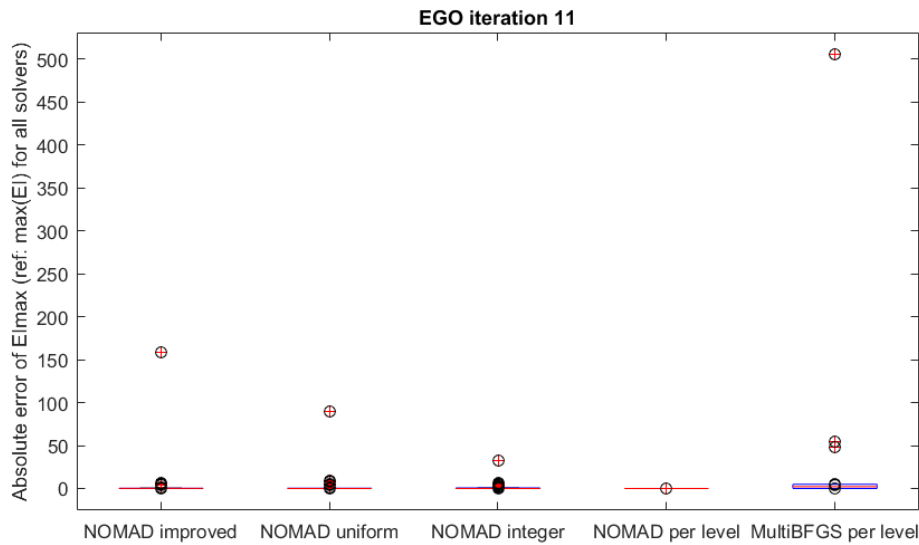


Figure 4. Absolute errors between the maximum EI of each method and the overall maximum EI. The results are presented in a Box-and-whisker plot accounting for 800 repetitions obtained by NOMAD with 3 different poll step strategies (adapted to integer, uniform sampling, "improved" sampling), one NOMAD per level and one multi-start BFGS per level.

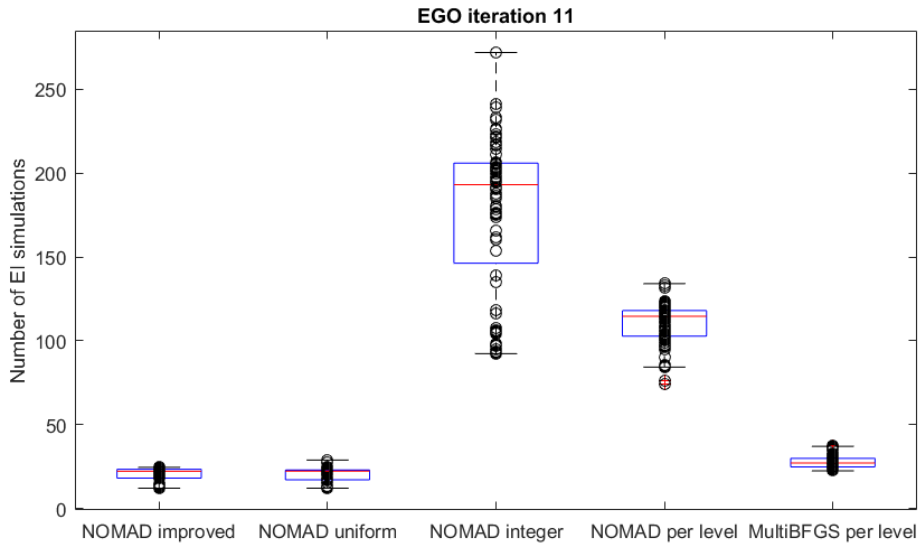


Figure 5. Number of simulations necessary to reach the maximal EI for 800 repetitions obtained by NOMAD with 3 different poll step strategies (adapted to integer, uniform sampling, "improved" sampling), one continuous NOMAD per level and one multi-start BFGS per level.

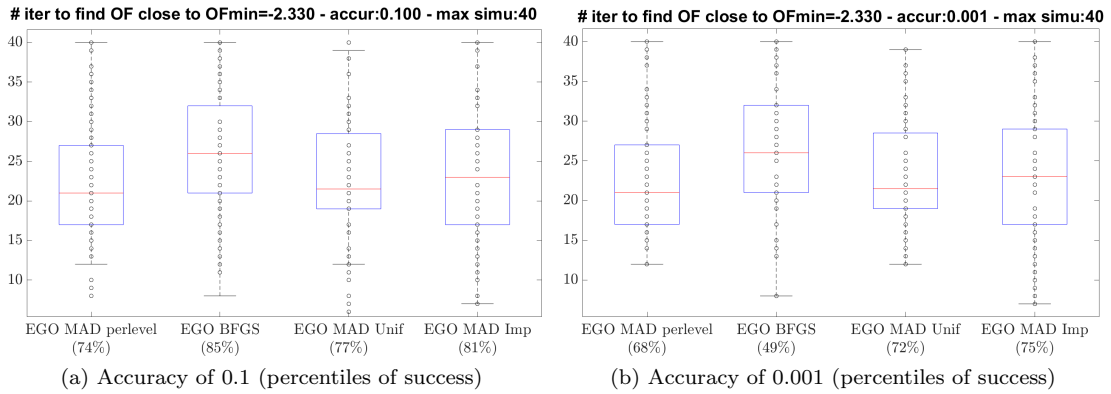


Figure 6. Number of simulations necessary to reach the maximal objective function for Cat-EGO with 2 given accuracies for a fixed simulation budgets of 40 simulations. Starting from 100 initial design of experiments of 5 points, 4 EI sub-optimization methods are evaluated: NOMAD with the 2 randomized poll step strategies (uniform sampling and "improved" sampling) and 2 "per level" strategies: Multi-start BFGS and NOMAD.

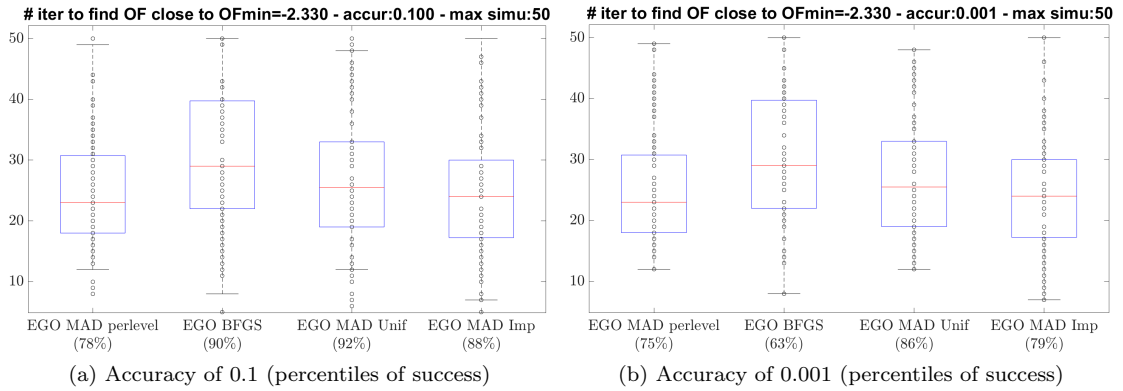


Figure 7. Number of simulations necessary to reach the maximal objective function for Cat-EGO with 2 given accuracies for a fixed simulation budgets of 50 simulations. Starting from 100 initial design of experiments of 5 points, 4 EI sub-optimization methods are evaluated: NOMAD with the 2 randomized poll step strategies (uniform sampling and "improved" sampling) and 2 "per level" strategies: Multi-start BFGS and NOMAD.

556 6.3. Benchmark: global skills of Cat-EGO

557 Note that in the sequel the NOMAD solver is run from an initial point which gives
 558 the highest EI value within the LHS sample used by Cat-EGO method (of size
 559 3 or $5 \times p \times N_{GL}$ as described further).

560
 561 We test our method on 15 box-constrained problems from the literature (Hock
 562 and Schittkowski (1981), Lukšan and Vlček (2000), and GECCO benchmark COCO
 563 (2017)) listed in Table 1. These test problems for continuous optimization are trans-
 564 formed into mixed integer problems. For Hock and Schittkowski (1981) and COCO
 565 (2017), following Liuzzi et al. (2012), we define integer variables z_i by restricting some
 566 continuous variables (every even index of variable vector) to take a finite number of
 567 values, m_i ; i.e.,

$$\forall \text{ even } i = 1, 2, \dots, q, x_i \in \left\{ \underline{x}_i + h \frac{(\bar{x}_i - \underline{x}_i)}{z_i - 1} \right\}, \text{ for } z_i = 0, 1, \dots, m_i - 1, \quad (20)$$

568 with \underline{x}_i and \bar{x}_i the respective lower and upper bounds of the original variable x_i .

569

Minimax problems

$$\min_{x \in [\underline{x}; \bar{x}]} F(x) := \max_{1 \leq z \leq m} (f_z(x)),$$

defined in Lukšan and Vlček (2000), are transformed into mixed categorical-continuous problems

$$\min_{x \in [\underline{x}; \bar{x}], 1 \leq z \leq m} \tilde{F}(x, z) := \begin{cases} f_1(x) & \text{if } z = 1 \\ f_2(x) & \text{if } z = 2 \\ \dots & \\ f_m(x) & \text{if } z = m \end{cases}$$

570 Functions from Hock and Schittkowski (1981) are smooth functions whereas the two

Test names	ncont (p)	ncat (q)	nlevels (m)	ref
EDV52	3	1	6	Lukšan and Vlček (2000)
RosenSuzuki	4	1	4	Lukšan and Vlček (2000)
SPIRAL	2	1	2	Lukšan and Vlček (2000)
Wong1	7	1	5	Lukšan and Vlček (2000)
HS2	1	1	4	Hock and Schittkowski (1981)
HS2 rand	1	1	4	Hock and Schittkowski (1981)
HS229log	1	1	4	Hock and Schittkowski (1981)
HS229log rand	1	1	4	Hock and Schittkowski (1981)
HS2	1	1	11	Hock and Schittkowski (1981)
HS229	1	1	11	Hock and Schittkowski (1981)
HS3	1	1	11	Hock and Schittkowski (1981)
bbob 10 3	2	1	4	COCO (2017)
bbob 21 3	2	1	4	COCO (2017)
bbob 22 3	2	1	4	COCO (2017)
bbob 21 5	3	2	4	COCO (2017)

Table 1. Benchmark functions: number of continuous variables, number of categorical variables and number of levels for each categorical variable.

571 other benchmark functions are more complex with several local minima (see, for instance, Figures 9 and 8).

573 The procedure to test the Cat-EGO method on these benchmark functions consists of

- 574 • an initial design of experiments built from concatenation of N_{GL} Latin hypercube designs, one for each global-level; the size of each design is $k \times p$, with $k = 3$ and 575 5, depending on the total number of levels,
- 576 • a limited budget of simulations, which is a common stopping criterion in practical applications of blackbox optimization for expensive simulators (Moré and Wild, 577 2009). Here, we chose a budget of 300 simulations.

580 Figure 10 illustrates the behavior of cat-EGO during the iterations: starting from an initial simulation set of 40 points built from concatenated Latin Hypercube designs of 10 points per level, the maximization of the expected improvement criterion leads to a compromise between space exploration and local minimization. This criterion 582 relies strongly on the learnt model and especially on the learnt correlations between global-levels: Figure 11 displays the evolution of the correlations $T_{c_i c_j}$ between the 583 4 levels during the iterations. A strong correlation has been detected between levels 584 1 and 4, between 2 and 3 and between 3 and 4 whereas levels 2 and 5 are anti-correlated. We observe then in Figure 10 that the exploration of the continuous domain 585 is complementary within the correlated levels whereas same zones of continuous space may be explored when levels are not correlated: *e.g.* levels 1 and 3. The adapted 586 structure of our model allows then to save some simulations thanks to the correlation information learnt from the simulated data, as shown on Table 2. In comparison, 587 RBFOpt and MISO methods explore much more at each level leading to a larger 588 number of simulations, as shown in Figures 12 and 13.

590 The results on the 15 functions are summarized in Table 2. Figures 14 and 15 display the mean relative error of objective functions (compared to the best value found by 591 the 4 optimizers) for the 15 functions of the benchmark. The Cat-EGO method is the most robust, leading to a mean error of 1%, whereas the other optimizers fail to obtain 592 593 594 595 596 597 598

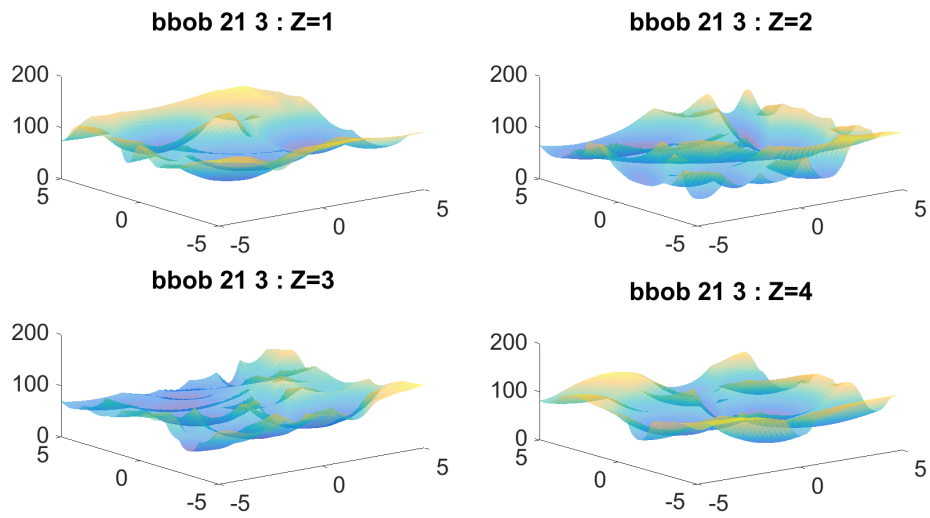


Figure 8. Function bbob21 of benchmark from [COCO \(2017\)](#) in 3 dimensions: 2 continuous variables and 1 categorical variable with 4 levels.

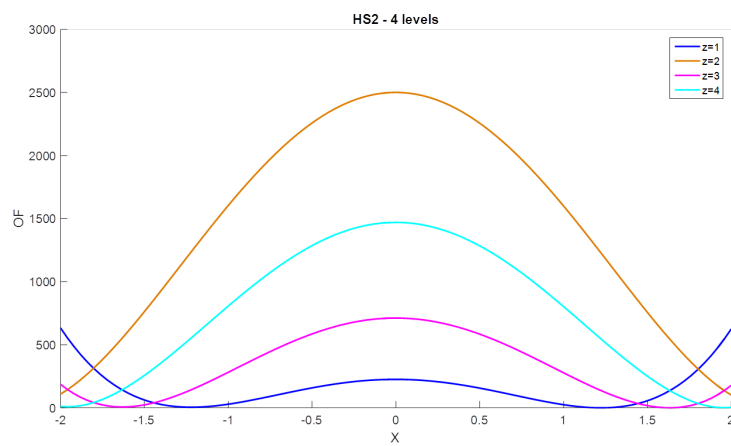


Figure 9. Function HS2 from [Hock and Schittkowski \(1981\)](#) in 2 dimensions: 1 continuous variable and 1 categorical variable with 4 levels.

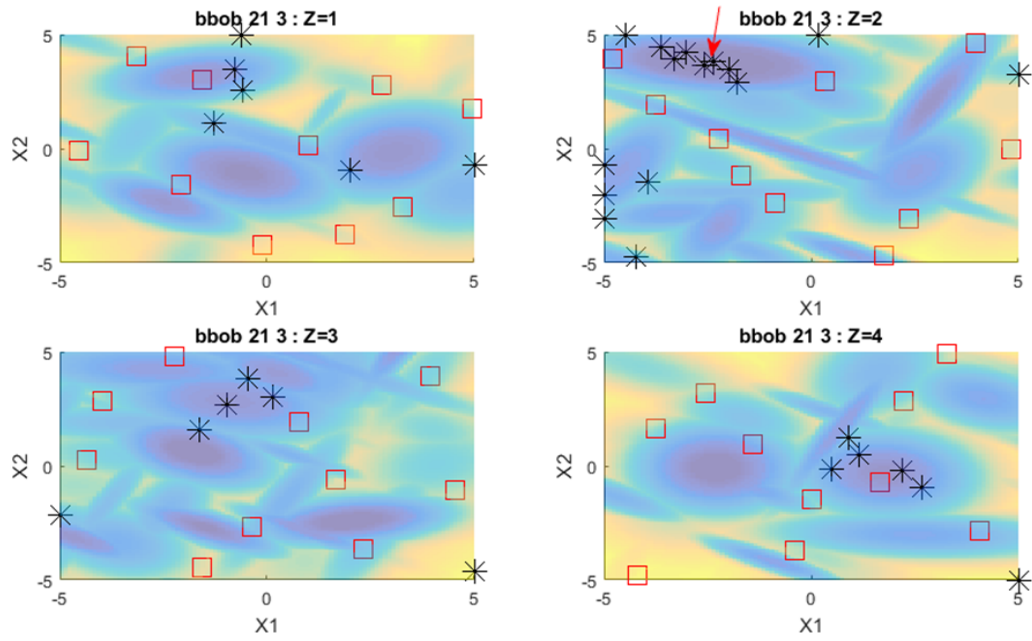


Figure 10. Simulation locations along EGO optimization iterations for test case bbob_21_3. The initial design of experiments is indicated with red squares. Black crosses are the additional simulated points determined by Expected Improvement maximization. The red arrow indicates the global optimum.

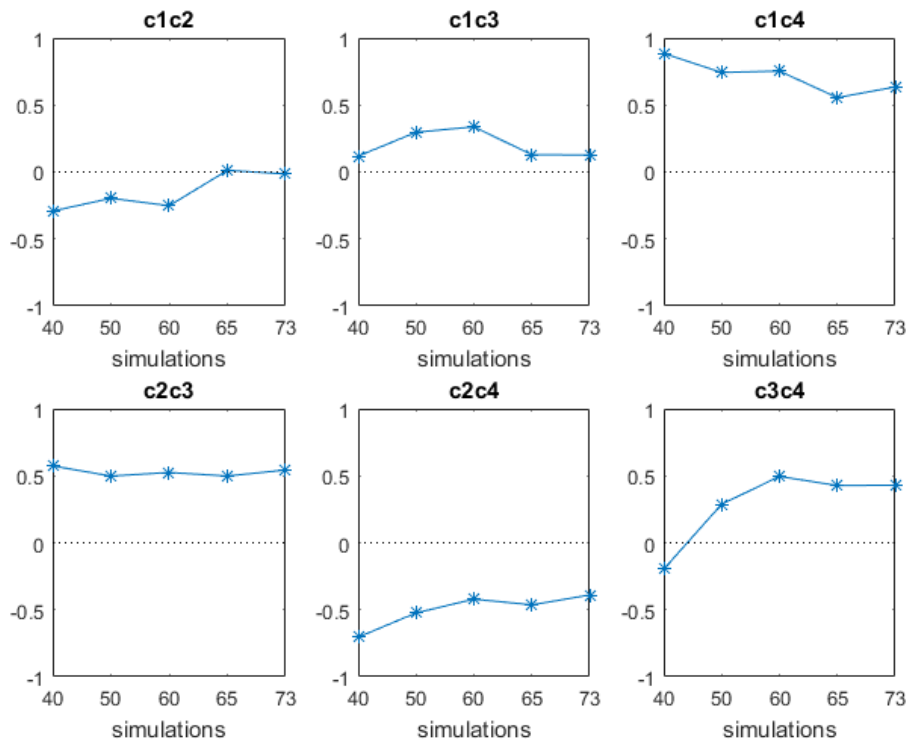


Figure 11. Evolution of EGO model correlations of categorical variables with simulations for test case bbob_21_3.

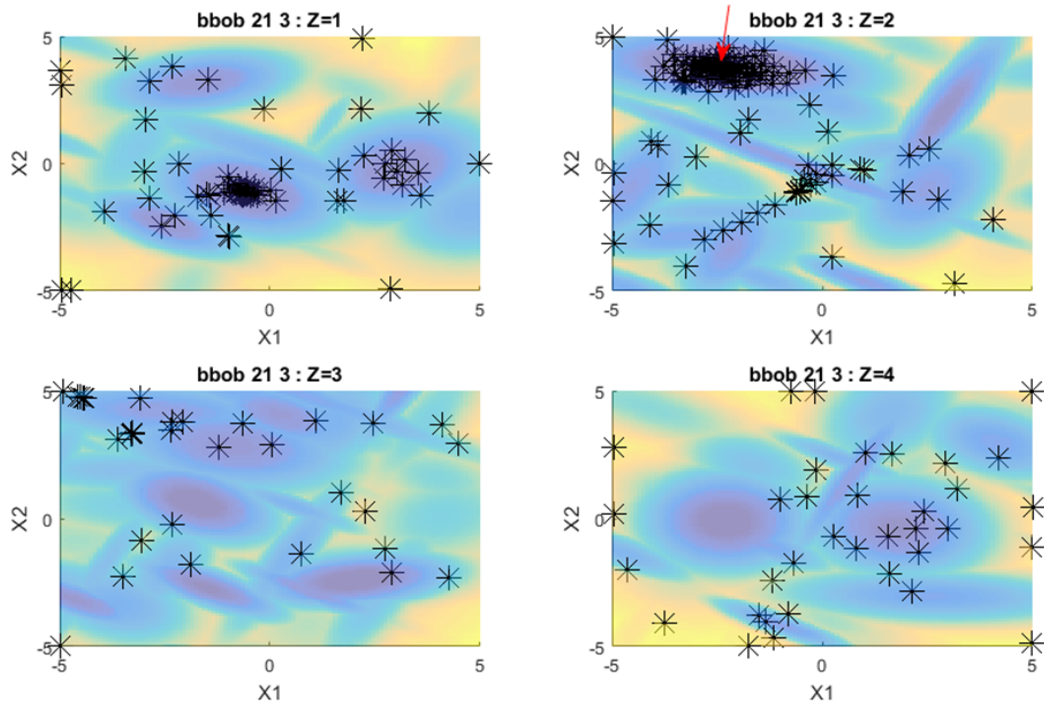


Figure 12. Simulation locations along RBFOpt optimization iterations for test case bbob_21_3. The red arrow indicates the global optimum.

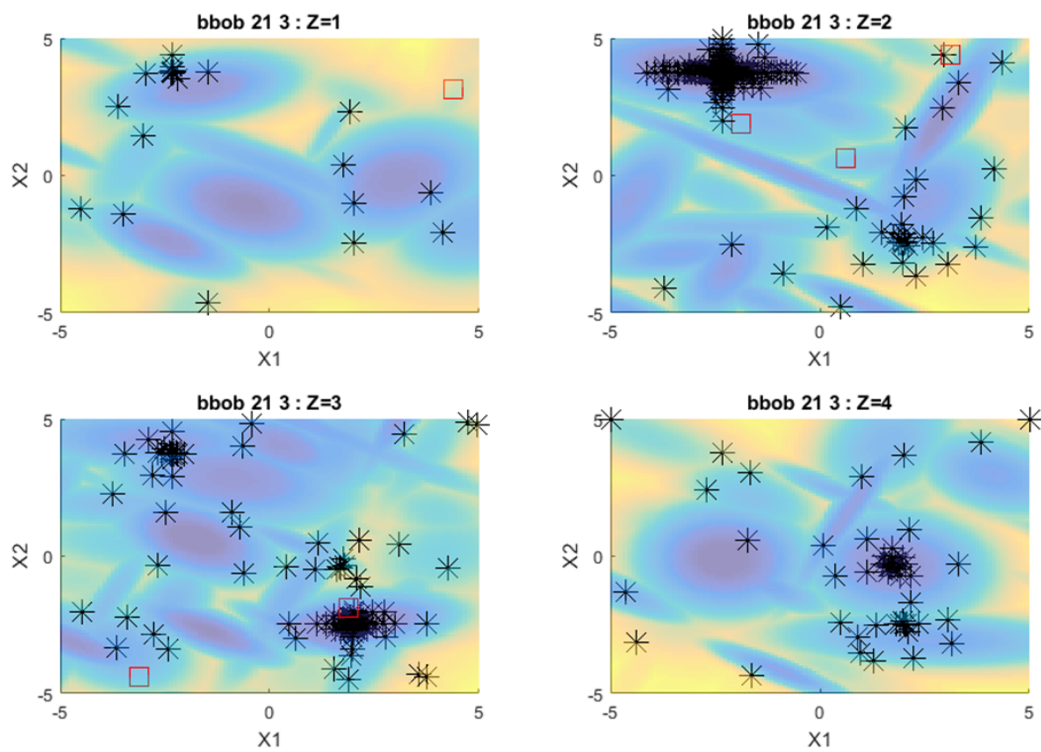


Figure 13. Simulation locations along MISO optimization iterations for test case bbob_21_3.

	Cat-EGO		NOMAD		RBFOpt		MISO	
EDV52	106	-1457.98	300	-99.00	166	-1455.51	300	-1458.00
RosenSuzuki	184	-113.13	300	-79.88	68	-79.87	300	-79.88
SPIRAL	235	0.00	300	0.00	297	0.00	300	0.03
Wong1	231	-1387.95	300	-3375.62	290	-2377.18	300	-3369.54
HS2 (10 levels)	124	0.05	177	6.16	104	0.05	242	0.05
HS2	38	0.05	244	0.92	217	0.40	220	0.05
HS2 rand	34	0.05	257	0.05	78	0.40	116	4.94
HS229 (10 levels)	52	0.00	211	0.09	215	0.00	167	0.00
HS229 log	56	-0.73	140	-0.75	86	-0.77	124	-0.75
HS229 log rand	35	-0.75	253	-0.77	63	-0.77	116	-0.77
HS3	45	0.00	15	0.00	296	0.00	11	0.00
bbob 10 3	151	-47.55	300	-54.62	194	-54.62	300	-54.62
bbob 21 3	68	40.78	169	42.75	221	40.78	300	40.78
bbob 22 3	154	-992.84	99	-995.06	296	-998.74	300	-995.06
bbob 21 5	164	41.01	300	44.61	177	49.07	300	42.61

Table 2. Comparison of Cat-EGO, NOMAD, RBFOpt and MISO optimizers on benchmark functions. Bold values indicate the runs which reach the best values in less than 300 simulations.

599 an acceptable accuracy in the allocated simulation budget (300) on several test cases,
600 as shown with the outliers on box-and-whisker plots of Figure 14.

601 **6.4. Behavior of the method for a larger size problem**

602 In this section we apply the 4 methods on a larger size test case with 2 continuous
603 variables and 4 categorical variables with 3 feasible values for each, that leads to 81
604 feasible categorical combinations. The test case is a modified version of "bbob 21"
605 function from GECCO benchmark COCO (2017) (see Table 1). The 4 categorical
606 variables are build from arbitrary discrete values of the 4 last original variables. The
607 global optimum is 0 at point $x_{1,\dots,6} = 1$ ($z_{3,\dots,6} = 2$).

608 Cat-EGO method is applied with two sizes of initial design of experiments: one of
609 162 points (a concatenated Latin Hypercube designs of 2 points per level) and one
610 global Latin hypercube design of 70 points (5 times the number of hyper-parameters).

611 Figure 16 illustrates the results obtained with cat-EGO, MISO, RBFOpt and NO-
612 MAD for a maximal budget of 600 simulations. The cat-EGO method obtains a smaller
613 value of the objective function in a smaller number of evaluations.

614 Figures 17, 18 illustrate the learning ability of cat-EGO models: the evolution of
615 the correlation matrix during the iterations is displayed. Starting from a larger design
616 of experiments (162 points) leads at the first iteration to a better estimate of the
617 correlations: Figure 18 illustrates that the correlation matrix at the first iteration
618 is very similar to the final one (after 600 simulations), whereas the correlations are
619 different between the first and the last iterations of cat-EGO run started from a
620 design of experiments of 70 points. When starting with 70 initial points, even with an
621 imperfect correlation matrix but informative enough, we observe that cat-EGO reaches
622 very quickly the optimal value of the objective function. The difference between the 70
623 points and 162 points initial DoE results (Figure 16) is probably due to an exploration-
624 exploitation trade-off but this is not a definitive conclusion since this is only one run
625 result.

626 in Figures 19 and 20 we observe how the function value evolves during the iterations

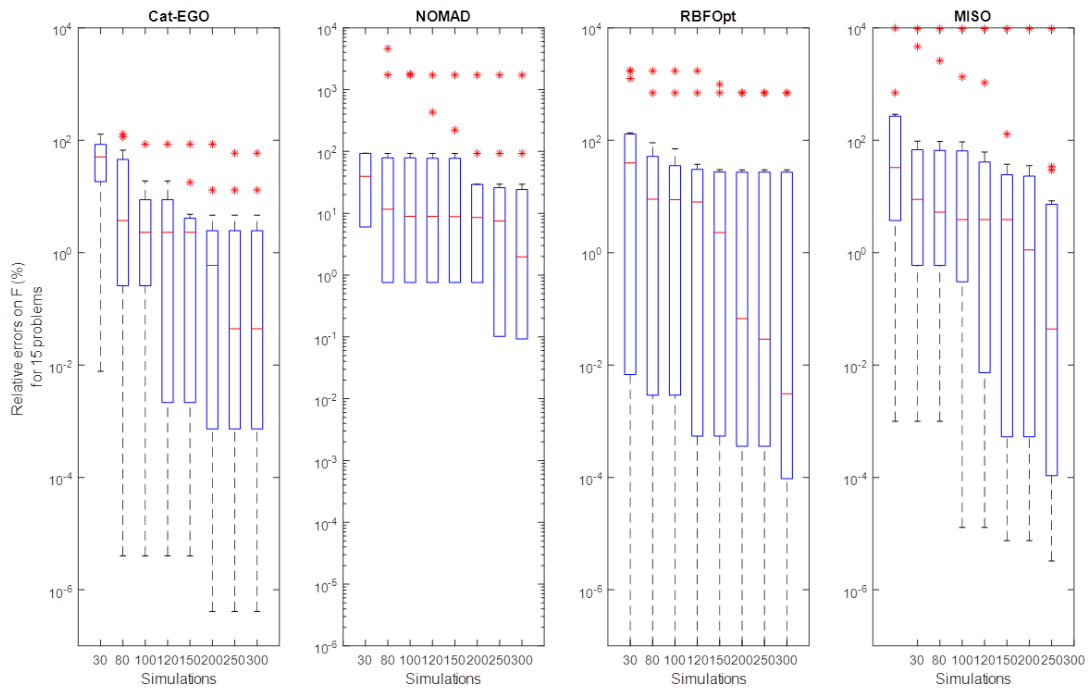


Figure 14. Box-and-whisker plots of mean relative error over the 15 benchmark functions versus simulations for Cat-EGO, NOMAD, RBFOpt and MISO optimizers.

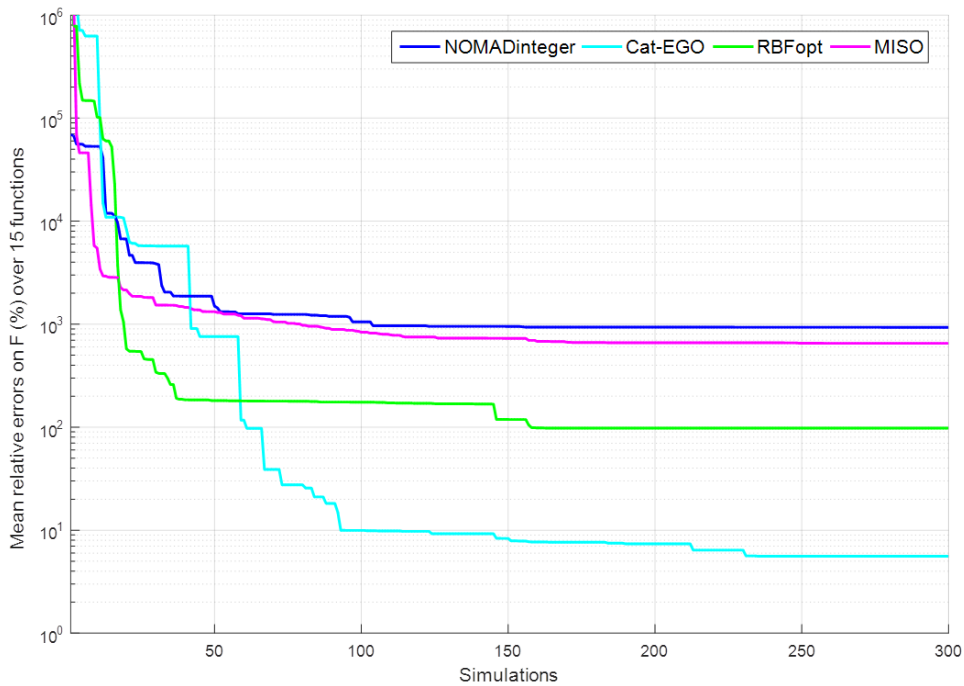


Figure 15. Mean relative error over the 15 benchmark functions versus simulations for Cat-EGO, NOMAD, RBFOpt and MISO optimizers.

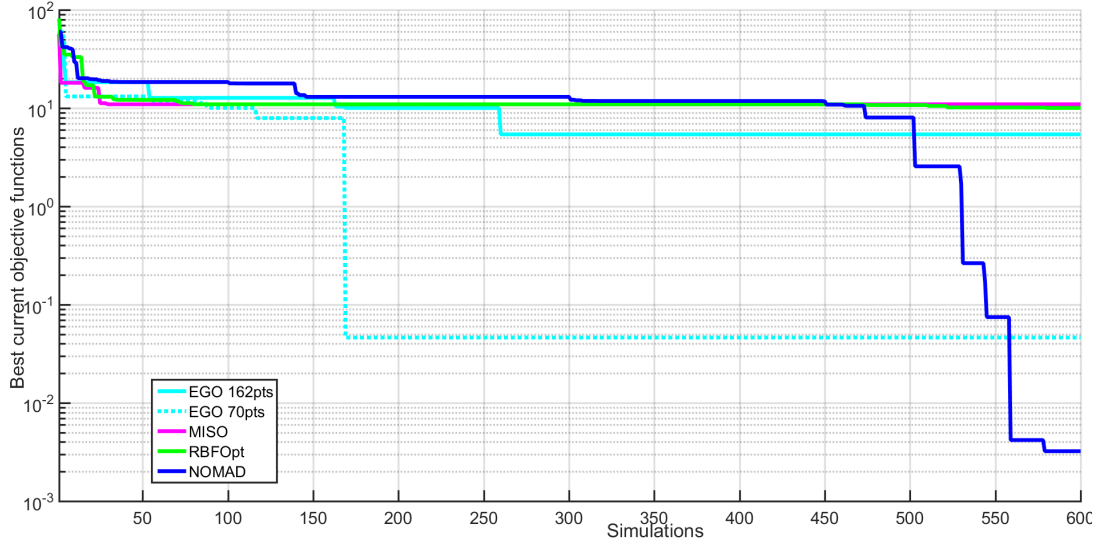


Figure 16. Evolution of the best current objective function value obtained with cat-EGO, MISO, RBFOpt and NOMAD for a maximal budget of 600 simulations.

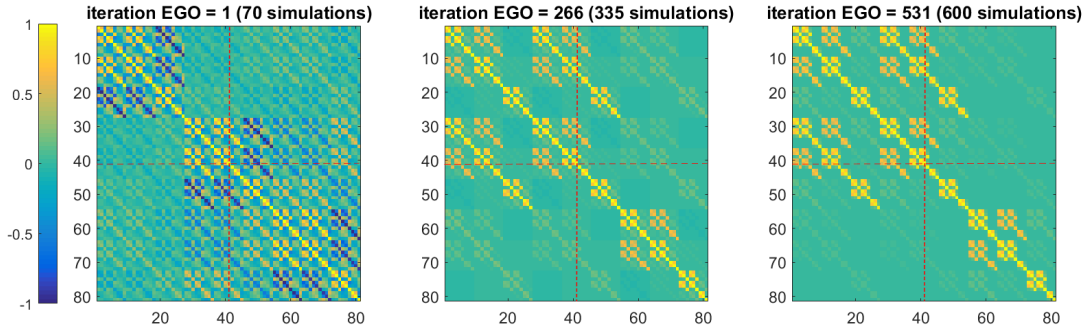


Figure 17. Model correlations learnt for categorical variables - between 81 combinations (the global-levels) taken pairwise - from the available simulation data points for 3 iterations of cat-EGO method started from a design of experiments of 70 points. The red lines indicate the row and the column associated with the global optimum.

627 with respect to the two different initial DoE sizes and when this value corresponds to
 628 points on the level containing the global minimum. The initial 70-points DoE has a
 629 very promising point on the optimal level which can explain the fast exploitation of
 630 this level. On the other hand, the initial 162-points DoE does not have much direct
 631 information on the optimal level but the correlations, being well estimated, push the
 632 exploitation of the optimal level quickly after the initial DoE. Indeed, the estimated
 633 correlations between the levels can drastically accelerate the exploration of the mixed
 634 variable space. Figures 21 and 22 illustrate how, at each iteration, the added point on
 635 a given global-level gives information on levels that are correlated with the currently
 636 explored one.

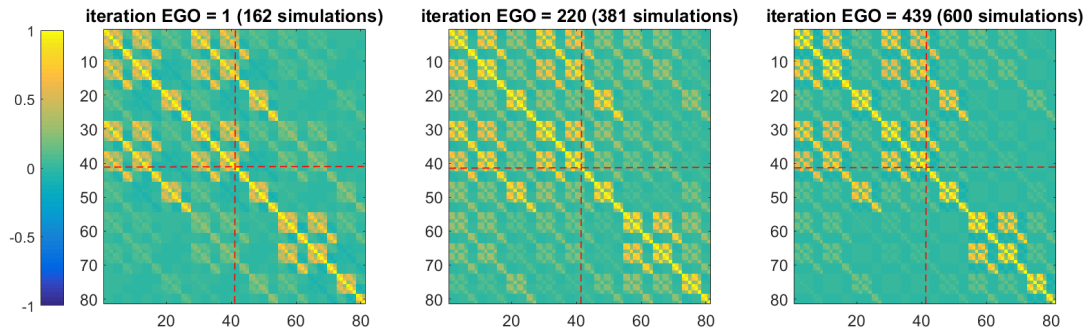


Figure 18. Model correlations learnt for categorical variables - between 81 combinations (the global-levels) taken pairwise - from the available simulation data points for 3 iterations of cat-EGO method started from a design of experiments of 162 points. The red lines indicate the row and the column associated with the global optimum.

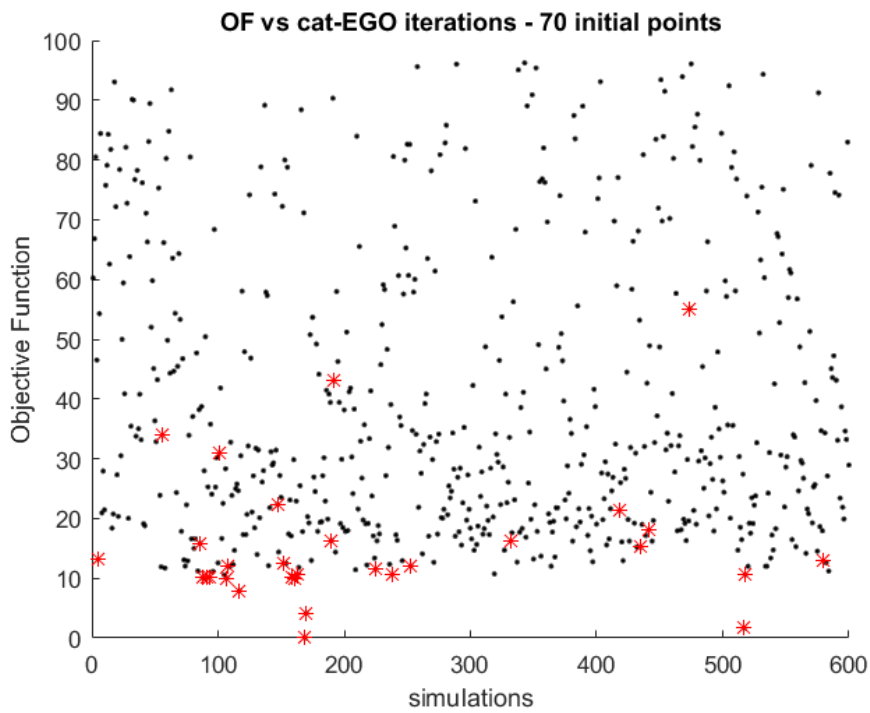


Figure 19. Objective function values during the iterations of Cat-EGO started with an initial design of 70 points. The red crosses highlight the values corresponding to a point on the global-level containing the minimum.

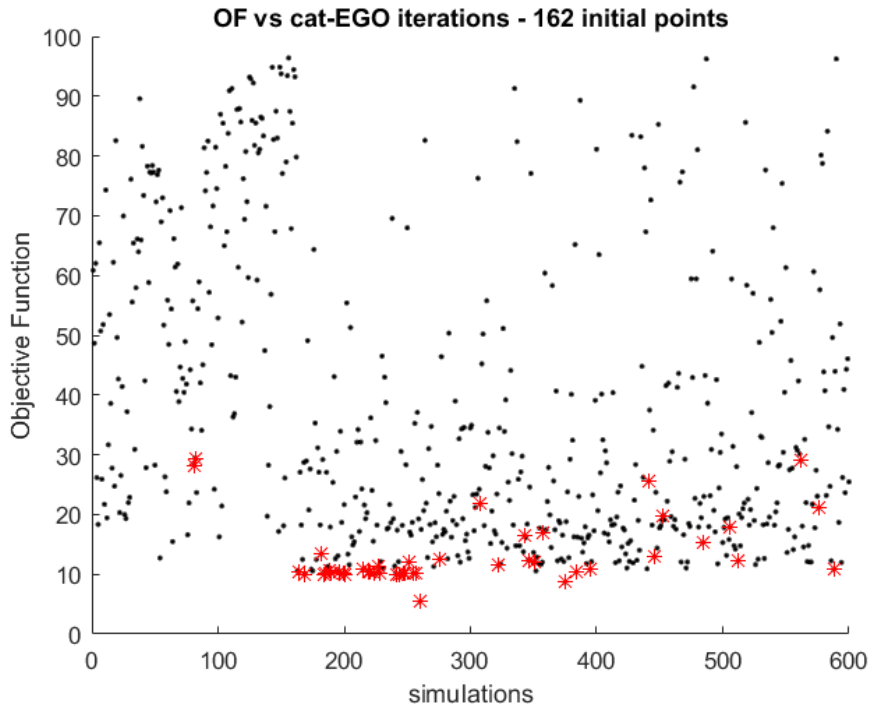


Figure 20. Objective function values during the iterations of Cat-EGO started with an initial design of 162 points. The red crosses highlight the values corresponding to a point on the global-level containing the minimum.

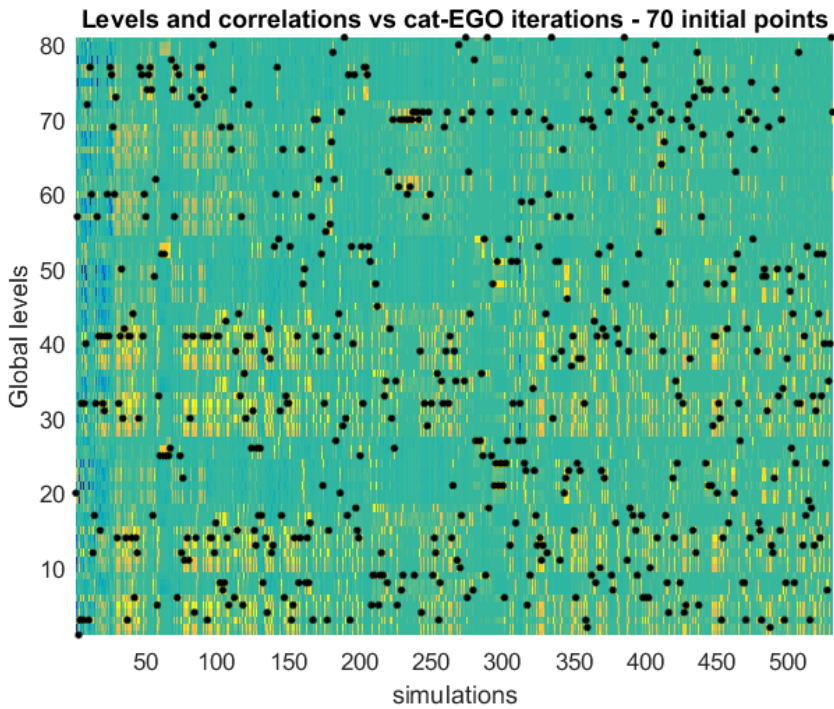


Figure 21. Global-levels visited during the iterations of cat-EGO with a 70 points initial DoE. The black dots indicate which global-level is visited at the corresponding iteration. At each iteration, the colors correspond to the degree of correlation between the visited global-level (where the black dot lies) and the other global-levels.

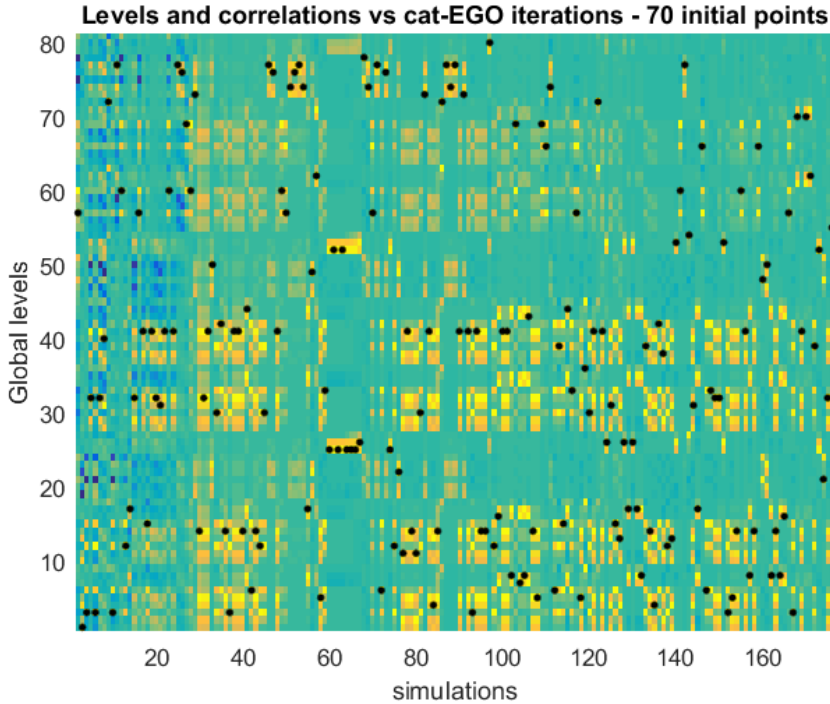


Figure 22. Zoom on 21 for the first simulations.

637 7. Conclusions and perspectives

638 We presented a strategy to tackle box-constrained mixed-integer global optimization
 639 problems involving expensive black box models and moderate number of input vari-
 640 ables. The proposed approach appears to be a robust method for finding global so-
 641 lutions of optimization problems. This success comes from the use of a probabilistic
 642 surrogate model (Gaussian processes) flexible enough (thanks to its hyper-parameter
 643 structure and their estimation) to capture relevant information on the optimized func-
 644 tion with respect to the continuous and the categorical variables. The GP approach
 645 also offers a quantification of the uncertainties on the objective function, enabling the
 646 construction of a built-in, optimization oriented, improvement criterion: the expected
 647 improvement (EI). For the EI sub-maximization task, we introduced a random explo-
 648 ration of the categorical variable space via a data based probability distribution. This
 649 latter enabled a faster recovery of an optimal solution of the main problem. On the
 650 other hand, the introduced flexibility comes with the price of a larger DoE as initial-
 651 ization and a larger computational time for the choice of the new points in the iterative
 652 scheme. Nevertheless, the method is still affordable when dealing with expensive-to-
 653 evaluate simulators and gives more robust results. In particular, we demonstrated the
 654 efficiency of the Cat-EGO strategy on a serie of test examples. We obtained on these
 655 test functions more robust results compared to RBFOpt and MISO. We are convinced
 656 that, on other test functions for which the kernel used by RBFOpt and MISO are well
 657 adapted, these latters will perform better in terms of minimal number of simulations
 658 required to reach the global optimum. This is explained because our method requires
 659 and uses a significant part of the simulations to learn the hyper-parameters. If the
 660 objective function is simple enough for the RBFOpt or MISO to be relevant then our

661 hyper-parameters learning stage is not necessary. In total generality, faced with an
 662 unknown objective function, we emphasize that it is more reasonable to tackle it with
 663 a flexible kernel (with more hyper-parameters), as the one we used.
 664 At this stage we suggest the following future research directions:

- 665 • On the one hand, reducing the number of categorical hyper-parameters with a
 666 learning strategy such as the one introduced in (Roustant et al., 2018) and, on
 667 the other hand, enabling different correlation lengths (one per global-level) in
 668 the continuous kernel part as in (Han et al., 2009),
- 669 • Further refinement of the discrete probability used in the NOMAD poll step in
 670 particular the adaptive calibration of the weight in (17),
- 671 • Penalizing the log-likelihood with the norm of the hyper-parameters with an
 672 adaptive penalization parameter driving the flexibility of the kernel. At the be-
 673 ginning of the method the size of the DoE is small with respect to the number of
 674 hyper-parameters. At this stage, the hyper-parameter optimization problem is
 675 not well posed (strongly non-convex). We thus propose to penalize their norm.
 676 This indeed inflates the hyper-parameters values, which leads to a very regular
 677 approximation of the function (with a small mean number of up-crossings for
 678 the continuous variables and high correlations for the categorical ones) and also
 679 a better posed hyper-parameter optimization problem.

680 Appendix A

We define the following rough global-level potential of improvement measure as

$$S_{k,i} = \bar{f}_{k,i} - 2\sigma_{k,i},$$

681 where $\bar{f}_{k,i}$ is the mean of the objective function values in the global-level c_i and $\sigma_{k,i}$
 682 the corresponding standard deviation. This measure $S_{k,i}$ takes into account the mean
 683 value of function evaluations in the global-level but also a measure of the variability of
 684 the continuous part within the global-level ($\sigma_{k,i}$). A small value of $S_{k,i}$ corresponds to
 685 a global level with high minimization potential. $S_{k,i}$ can be seen as an approximation
 686 of $\bar{M}(z)$ in (16) with $z = c_i$. We then calculate the quantity

$$S_{k,i}^R = \frac{1}{1 + \exp(-b_{k,i} \frac{f_{min} - S_{k,i}}{\hat{\sigma}_{k,i}})} \quad (21)$$

687 which is a sigmoid function parametrized by $b_{k,i}$ evaluated at $(f_{min} - S_{k,i})/\hat{\sigma}_{k,i}$. This
 688 quantity approximates $\Psi_z(\frac{f_{min} - \bar{M}(z)}{\sigma_M(z)})$ in (16) for $z = c_i$. The coefficient $b_{k,i}$ should be
 689 selected in order to approximate the cumulative distribution Ψ_{c_i} and $\hat{\sigma}_{k,i}$ should be
 690 an approximation of $\sigma_M(c_i)$.

691 As presented, the proposed exploration scheme depends on the parameters α_k ,
 692 $b_{k,i}$ and $\hat{\sigma}_{k,i}$. For the numerical results we directly set $\alpha_k = 0$ so that only $p_{k,i}^m$, the
 693 probability that the global-level c_i has high potential of containing the minimum,
 694 is considered in (17). We imposed $b_{k,i} = 1$ for all k and i since estimating the
 695 distribution seems too expensive. We also imposed $\hat{\sigma}_{k,i} = 1$ for all k and i , since an
 696 accurate approximation of this term is also expensive to compute and by definition
 697 $S_{k,i}$ already integrates some insight on the standard deviation of the corresponding
 698 global-level. Nevertheless, setting $\hat{\sigma}_{k,i}$ constant implies that the probability $p_{k,i}^m$ will

699 not converge to 1 for the level containing the minimum and 0 otherwise. This is not
700 an issue here since the limit in k will not be reached within the limited number of
701 iterations to be expected (a few hundreds), and furthermore, the probability will
702 still be large in the global-level containing the objective function minimum. A finer
703 analysis of the tuning of the parameters α_k , $b_{k,i}$ and $\hat{\sigma}_{k,i}$ is postponed to further work.
704

705 References

- 706 Abramson M, Audet C, Chrissis J, Walston J. 2009. Mesh adaptive direct search
707 algorithms for mixed variable optimization. *Optimization Letters*. 3(1):35. Available
708 from: <https://doi.org/10.1007/s11590-008-0089-2>.
- 709 Adler RJ. 1981. *The geometry of random fields*. Wiley, Chichester.
- 710 Audet C, Dennis JE. 2000. Pattern search algorithms for mixed variable programming.
711 *SIAM Journal on Optimization*. 11:573–594.
- 712 Audet C, J E Dennis J. 2006. Mesh adaptive direct search algorithms for con-
713 strained optimization. *SIAM Journal on Optimization*. 17(1):188–217. Available
714 from: <https://doi.org/10.1137/040603371>.
- 715 Byrd RH, Lu P, Nocedal J, Zhu C. 1995. A limited memory algorithm for bound
716 constrained optimization. *SIAM J Scientific Computing*. 16:1190–1208.
- 717 COCO. 2017. Black-box optimization benchmarking. Available from: <http://coco.gforge.inria.fr/>.
- 718 Comola F, Janna C, Lovison A, Minini M, Tamburini A, Teatini P. 2016. Efficient
719 global optimization of reservoir geomechanical parameters based on synthetic aper-
720 ture radar-derived ground displacements. *GEOPHYSICS*. 81(3):M23–M33. Avail-
721 able from: <https://doi.org/10.1190/geo2015-0402.1>.
- 722 Costa A, Nannicini G. 2018. RBFOpt: an open-source library for black-box optimiza-
723 tion with costly function evaluations. *Mathematical Programming Computation*.
724 10(4):597–629. Available from: <https://doi.org/10.1007/s12532-018-0144-7>.
- 725 Gramacy RB, Taddy M. 2010. Categorical inputs, sensitivity analysis, optimization
726 and importance tempering with `tgp` version 2, an r package for treed gaussian process
727 models. *Journal of Statistical Software*. 33(6).
- 728 Gutmann HM. 2001. A radial basis function method for global optimization. *J*
729 *of Global Optimization*. 19(3):201–227. Available from: [http://dx.doi.org/10.](http://dx.doi.org/10.1023/A:1011255519438)
730 [1023/A:1011255519438](http://dx.doi.org/10.1023/A:1011255519438).
- 731 Hamza K, Shalaby M. 2014. A framework for parallelized efficient global optimiza-
732 tion with application to vehicle crashworthiness optimization. *Engineering Opti-*
733 *mization*. 46(9):1200–1221. Available from: [https://doi.org/10.1080/0305215X.](https://doi.org/10.1080/0305215X.2013.827672)
734 [2013.827672](https://doi.org/10.1080/0305215X.2013.827672).
- 735 Han G, Santner T, Notz W, Bartel D. 2009. Prediction for computer experiments
736 having quantitative and qualitative input variables. *Technometrics*. 51(3):278–288.
737 Available from: <https://doi.org/10.1198/tech.2009.07132>.
- 738 Helbert C, Dupuy D, Carraro L. 2009. Assessment of uncertainty in computer exper-
739 iments: From universal kriging to bayesian kriging. *Applied Stochastic Models in*
740 *Business and Industry*. 25:99–113.
- 741 Hock W, Schittkowski K. 1981. *Test examples for nonlinear programming codes*. vol.
742 187. *Lecture Notes in Economics and Mathematical Systems*, Berlin: Springer.
- 743 Jones DR, Schonlau M, Welch W. 1998. Efficient global optimization of expensive
744 black-box functions. *Journal of Global Optimization*. 13(4):455–492. Available from:
745

- 746 <https://doi.org/10.1023/A:1008306431147>.
- 747 Kanazaki M, Matsuno T, Maeda K, Kawazoe H. 2015. Efficient global optimization ap-
748 plied to wind tunnel evaluation-based optimization for improvement of flow control
749 by plasma actuators. *Engineering Optimization*. 47(9):1226–1242. Available from:
750 <https://doi.org/10.1080/0305215X.2014.958733>.
- 751 Liuzzi G, Lucidi S, Rinaldi F. 2012. Derivative-free methods for bound constrained
752 mixed-integer optimization. *Computational Optimization and Applications*.
753 53(2):505–526. Available from: <https://doi.org/10.1007/s10589-011-9405-3>.
- 754 Lukšan L, Vlček J. 2000. Test problems for nonsmooth unconstrained and linearly con-
755 strained optimization. Technical report VT798-00, Institute of Computer Science,
756 Academy of Sciences of the Czech Republic.
- 757 McKay MD, Beckman RJ, Conover WJ. 1979. A comparison of three methods for
758 selecting values of input variables in the analysis of output from a computer code.
759 *Technometrics*. 21:239–245.
- 760 Moré JJ, Wild SM. 2009. Benchmarking derivative-free optimization algorithms. *SIAM*
761 *J on Optimization*. 20(1):172–191. Available from: [http://dx.doi.org/10.1137/](http://dx.doi.org/10.1137/080724083)
762 [080724083](http://dx.doi.org/10.1137/080724083).
- 763 Muller J. 2016. MISO: mixed-integer surrogate optimization framework. *Optimiza-*
764 *tion and Engineering*. 17(1):177–203. Available from: [https://doi.org/10.1007/](https://doi.org/10.1007/s11081-015-9281-1)
765 [s11081-015-9281-1](https://doi.org/10.1007/s11081-015-9281-1).
- 766 Muller J, Shoemaker CA, Piché R. 2013. SO-MI: A surrogate model algorithm for
767 computationally expensive nonlinear mixed-integer black-box global optimization
768 problems. *Computers & Operations Research*. 40(5):1383 – 1400. Available from:
769 <http://www.sciencedirect.com/science/article/pii/S0305054812001967>.
- 770 Pelamatti J, Brevault L, Balesdent M, Talbi EG, Guerin T. 2018. Efficient global opti-
771 mization of constrained mixed variable problems. *Journal of Global Optimization*.
772 25. Available from: <https://doi.org/10.1007/s10898-018-0715-1>.
- 773 Pinheiro J, Bates D. 1996. Unconstrained parametrizations for variance-covariance
774 matrices. *Statistics and Computing*. 6(3):289–296.
- 775 Pinheiro J, Bates D. 2009. *Mixed-effects models in s and s-plus*. Statistics and Com-
776 puting Springer New York.
- 777 Potdar K, Pardawala TS, Pai CD. 2017. A comparative study of categorical variable en-
778 coding techniques for neural network classifiers. *International Journal of Computer*
779 *Applications*. 175(4):7–9. Available from: [http://www.ijcaonline.org/archives/](http://www.ijcaonline.org/archives/volume175/number4/28474-2017915495)
780 [volume175/number4/28474-2017915495](http://www.ijcaonline.org/archives/volume175/number4/28474-2017915495).
- 781 Qian P. 2012. Sliced latin hypercube designs. *Journal of the American Statistical Assoc-*
782 *iation*. 107(497):393–399. Available from: [https://doi.org/10.1080/01621459.](https://doi.org/10.1080/01621459.2011.644132)
783 [2011.644132](https://doi.org/10.1080/01621459.2011.644132).
- 784 Qian P, Wu H, Wu J. 2008. Gaussian process models for computer experiments with
785 qualitative and quantitative factors. *Annals of Statistics*. *Technometrics*, Vol. 50,
786 No.3:383–396.
- 787 Rasmussen CE. 2006. *Gaussian processes for machine learning*. MIT Press.
- 788 Roustant O, Ginsbourger D, Deville Y. 2012. Dicekriging, diceoptim: Two r pack-
789 ages for the analysis of computer experiments by kriging-based metamodeling and
790 optimization. *Journal of Statistical Software*. 51(1).
- 791 Roustant O, Padonou E, Deville Y, Clément A, Perrin G, Giorla J, Wynn H. 2018.
792 Group kernels for Gaussian process metamodels with categorical inputs. ArXiv e-
793 prints.
- 794 Sacks J, Welch W, Mitchell T, Wynn H. 1989. Design and analysis of computer ex-

- 795 periments. *Statistical Science*. 4(4):409 – 435.
- 796 Santner T, Williams B, Notz W. 2003. *The design and analysis of computer experi-*
797 *ments*. New York: Springer.
- 798 Schonlau M. 1997. *Computer experiments and global optimization* [dissertation]. Uni-
799 *versity of Waterloo*.
- 800 Swiler L, Hough P, Qian P, Xu X, Storlie C, Lee H. 2014. *Surrogate models for mixed*
801 *discrete-continuous variables*. Cham: Springer International Publishing. p. 181–202.
- 802 Taddy M, Lee HKH, Gray GA, Griffin JD. 2009. Bayesian guided pattern search for
803 robust local optimization. *Technometrics*. 51:389–401.
- 804 Zhang Y, Notz W. 2015. Computer experiments with qualitative and quantitative
805 variables: A review and reexamination. *Quality Engineering*. 27(1):2–13. Available
806 from: <https://doi.org/10.1080/08982112.2015.968039>.
- 807 Zhou Q, Qian PZG, Zhou S. 2011. A simple approach to emulation for computer models
808 with qualitative and quantitative factors. *Technometrics*. 53(3):266–273. Available
809 from: <https://doi.org/10.1198/TECH.2011.10025>.