



**HAL**  
open science

# Differentiating the multipoint Expected Improvement for optimal batch design

Sébastien Marmin, Clément Chevalier, David Ginsbourger

► **To cite this version:**

Sébastien Marmin, Clément Chevalier, David Ginsbourger. Differentiating the multipoint Expected Improvement for optimal batch design. 2015. hal-01133220v1

**HAL Id: hal-01133220**

**<https://hal.science/hal-01133220v1>**

Preprint submitted on 18 Mar 2015 (v1), last revised 27 Aug 2019 (v4)

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

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

# Differentiating the multipoint Expected Improvement for optimal batch design

Sébastien Marmin<sup>1,2</sup>, Clément Chevalier<sup>3</sup>, David Ginsbourger<sup>1</sup>

<sup>1</sup> Department of Mathematics and Statistics,  
University of Bern, Switzerland

<sup>2</sup> Institut de Radioprotection et de Sûreté Nucléaire,  
Cadarache, France

<sup>3</sup> Institute of Mathematics,  
University of Zurich, Switzerland

March 18, 2015

## Abstract

This work deals with parallel optimization of expensive objective functions which are modelled as sample realizations of Gaussian processes. The study is formalized as a Bayesian optimization problem, or continuous multi-armed bandit problem, where a batch of  $q > 0$  arms is pulled in parallel at each iteration. Several algorithms have been developed for choosing batches by trading off exploitation and exploration. As of today, the maximum Expected Improvement (EI) and Upper Confidence Bound (UCB) selection rules appear as the most prominent approaches for batch selection. Here, we build upon recent work on the multipoint Expected Improvement criterion, for which an analytic expansion relying on Tallis' formula was recently established. The computational burden of this selection rule being still an issue in application, we derive a closed-form expression for the gradient of the multipoint Expected Improvement, which aims at facilitating its maximization using gradient-based ascent algorithms. Substantial computational savings are shown in application. In addition, our algorithms are tested numerically and compared to state-of-the-art UCB-based batch-sequential algorithms. Combining starting designs relying on UCB with gradient-based EI local optimization finally appears as a sound option for batch design in distributed Gaussian Process optimization.

**Keywords :** Bayesian Optimization, Batch-sequential design, GP, UCB.

# 1 Introduction

Global optimization of deterministic functions under a drastically limited evaluation budget is a topic of growing interest with important industrial applications. Dealing with such expensive black-box simulators is typically addressed through the introduction of surrogate models that are used both for reconstructing the objective function and guiding parsimonious evaluation strategies. This approach is used in various scientific communities and referred to as Bayesian optimization, but also as kriging-based or multi-armed bandit optimization [20, 5, 23] [16, 15, 25, 11]. Among such Gaussian process optimization methods, two concepts of algorithm relying on sequential maximization of in-fill sampling criteria are particularly popular in the literature. In the EGO algorithm of [16], the sequence of decisions (of where to evaluate the objective function at each iteration) is guided by the Expected Improvement (EI) criterion [19], which is known to be one-step lookahead optimal [14]. On the other hand, the Upper Confidence Bound (UCB) algorithm [1] maximizes sequentially a well-chosen kriging quantile, that is, a quantile of the pointwise posterior Gaussian process distribution. Similarly to EI [24, 6], the consistency of the algorithm has been established and rates of convergence have been obtained [23].

Recently, different methods inspired from the two latter algorithms have been proposed to deal with the typical case where  $q > 1$  CPUs are available. Such synchronous distributed methods provide at each iteration a batch of  $q$  points which can be evaluated in parallel. For instance, [10] generalizes the UCB algorithm to a batch-sequential version by maximizing kriging quantiles and assuming dummy responses equal to the posterior mean of the Gaussian process. This approach can be compared with the so-called Kriging Believer strategy of [15] where each batch is obtained by sequentially maximizing the one-point EI under the assumption that the previously chosen points have a response equal to their Kriging mean. Originally, the strategies suggested in [15] were introduced to cope with the difficulty to evaluate and maximize the multipoint Expected Improvement ( $q$ -EI) [22], which is the generalization of EI known to be one-batch lookahead optimal [7, 14]. One of the bottlenecks for  $q$ -EI maximization was that it was until recently evaluated through Monte-Carlo simulations [15], a reason that motivated [11] to propose a stochastic gradient algorithm for its maximization. Now, [8] established a closed-form expression enabling to compute  $q$ -EI at any batch of  $q$  points without appealing to Monte-Carlo simulations. However, the computational complexity involved to compute the criterion is still high and quickly grows with  $q$ . Besides, little has been published about the difficult maximization of the  $q$ -EI itself, which is an optimization problem in dimension  $qd$ , where  $d$  is the number of input variables.

In this work, we contribute to the latter problem by giving an analytical gradient of  $q$ -EI, in the space of dimension  $qd$ . Such a gradient is meant to simplify the local maximization of  $q$ -EI using gradient-based ascent algorithms. Closed-form expressions of  $q$ -EI and its gradient have been implemented in the DiceOptim R package [21], together with a multistart BFGS algorithm for max-

imizing  $q$ -EI. In addition, we suggest to use results of the BUCB algorithm as initial batches in multistart gradient-based ascents. These starting batches are shown to yield good local optima for  $q$ -EI. This article is organized as follows. Section 2 quickly recalls the basics of Gaussian process modeling and the closed-form expression of  $q$ -EI obtained in [8]. Section 3 details the analytical  $q$ -EI gradient. Finally, numerical experiments comparing the performances of the  $q$ -EI maximization-based strategy and the BUCB algorithms are provided and discussed in Section 4. For readability and conciseness, the most technical details about  $q$ -EI gradient calculation are sent in Appendix.

## 2 General Context

Let  $f : \mathbf{x} \in D \subset \mathbb{R}^d \rightarrow \mathbb{R}$  be a real-valued function defined on a compact subset  $D$  of  $\mathbb{R}^d$ ,  $d \geq 1$ . Throughout this article, we assume that we dispose of a set of  $n$  evaluations of  $f$ ,  $\mathcal{A}_n = (\mathbf{x}_{1:n} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \mathbf{y}_{1:n} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^\top)$ , and that our goal is to evaluate  $f$  at well-chosen batches of  $q$  points in order to globally maximize it. Following each batch of evaluations, we observe  $q$  deterministic scalar responses, or rewards,  $y_{n+1} = f(\mathbf{x}_{n+1}), \dots, y_{n+q} = f(\mathbf{x}_{n+q})$ . We use past observations in order to carefully choose the next  $q$  observation locations, aiming in the end to minimize the one-step lookahead regret  $f(\mathbf{x}^*) - t_{n+q}$ , where  $\mathbf{x}^*$  is a maximizer of  $f$  and  $t_i = \max_{j=1, \dots, i} (f(\mathbf{x}_j))$ . In this section, we first define the Gaussian process (GP) surrogate model used to make the decisions. Then we introduce the  $q$ -EI which is the optimal one-batch lookahead criterion (see, e.g., [14, 3, 12] for a definition and [14, 7] for a proof).

### 2.1 Gaussian process modeling

The objective function  $f$  is a priori assumed to be a sample from a Gaussian process  $Y \sim \mathcal{GP}(\mu, C)$ , where  $\mu(\cdot)$  and  $C(\cdot, \cdot)$  are respectively the mean and covariance function of  $Y$ . At fixed  $\mu(\cdot)$  and  $C(\cdot, \cdot)$ , conditioning  $Y$  on the set of observations  $\mathcal{A}_n$  yields a GP posterior  $Y(\mathbf{x}) | \mathcal{A}_n \sim \mathcal{GP}(\mu_n, C_n)$  with:

$$\mu_n(\mathbf{x}) = \mu(\mathbf{x}) + \mathbf{c}_n(\mathbf{x})^\top \mathbf{C}_n^{-1} (\mathbf{y}_{1:n} - \mu(\mathbf{x}_{1:n})), \text{ and} \quad (1)$$

$$C_n(\mathbf{x}, \mathbf{x}') = C(\mathbf{x}, \mathbf{x}') - \mathbf{c}_n(\mathbf{x})^\top \mathbf{C}_n^{-1} \mathbf{c}_n(\mathbf{x}'), \quad (2)$$

where  $\mathbf{c}_n(\mathbf{x}) = (C(\mathbf{x}, \mathbf{x}_i))_{1 \leq i \leq n}$ , and  $\mathbf{C}_n = (C(\mathbf{x}_i, \mathbf{x}_j))_{1 \leq i, j \leq n}$ . Note that, in realistic application settings, the mean and the covariance  $\mu$  and  $C$  of the prior are assumed to depend on several parameters which require to be estimated. The results presented in this article and their implementations in the R package `DiceOptim` are compatible with this more general case. More detail about Equations 1, 2 with or without trend and covariance parameter estimation can be found in [21] and is omitted here for conciseness.

## 2.2 The Multipoint Expected Improvement criterion

The Multipoint Expected Improvement ( $q$ -EI) selection rule consists in maximizing, over all possible batches of  $q$  points, the following criterion, which depends on a batch  $\mathbf{X} = (\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+q}) \in D^q$ :

$$\text{EI}(\mathbf{X}) = \mathbb{E}[(\max Y(\mathbf{X}) - T_n)_+ | \mathcal{A}_n], \quad (3)$$

where  $(\cdot)_+ = \max(\cdot, 0)$ , and the threshold  $T_n$  is the currently observed maximum of  $Y$ , i.e.  $T_n = \max_{1 \leq j \leq n} Y(\mathbf{x}_j)$ . Recalling that  $Y(\mathbf{X}) | \mathcal{A}_n \sim \mathcal{N}(\boldsymbol{\mu}_n(\mathbf{X}), C_n(\mathbf{X}, \mathbf{X}))$ , and denoting  $Y(\mathbf{X}) = (Y_1, \dots, Y_q)^\top$ , an analytic expression of  $q$ -EI at locations  $\mathbf{X}$  over any threshold  $T \in \mathbb{R}$  can be found in [8] and is reproduced here :

$$\text{EI}(\mathbf{X}) = \sum_{k=1}^q \left( (m_k - T) \Phi_{q, \Sigma^{(k)}}(-\mathbf{m}^{(k)}) + \sum_{i=1}^q \Sigma_{ik}^{(k)} \varphi_{\Sigma_{ii}^{(k)}}(m_i^{(k)}) \Phi_{q-1, \Sigma_{|i}^{(k)}}(-\mathbf{m}_{|i}^{(k)}) \right) \quad (4)$$

where  $\varphi_{\sigma^2}(\cdot)$  and  $\Phi_{p, \Gamma}(\cdot)$  are respectively the density function of the centered normal distribution with variance  $\sigma^2$  and the  $p$ -variate cumulative distribution function (CDF) of the centered normal distribution with covariance  $\Gamma$  ;  $\mathbf{m} = \mathbb{E}(Y(\mathbf{X}) | \mathcal{A}_n)$  and  $\Sigma = \text{cov}(Y(\mathbf{X}) | \mathcal{A}_n)$  are the conditional mean vector and covariance matrix of  $Y(\mathbf{X})$  ;  $\mathbf{m}^{(k)}$  and  $\Sigma^{(k)}$ ,  $1 \leq k \leq q$ , are the conditional mean vector and covariance matrix of the affine transformation of  $Y(\mathbf{X})$ ,  $\mathbf{Z}^{(k)} = L^{(k)}Y(\mathbf{X}) + \mathbf{b}^{(k)}$ , defined as  $Z_j^{(k)} := Y_j$  for  $j \neq k$  and  $Z_k^{(k)} := T - Y_k$  ; and finally, for  $(k, i) \in \{1, \dots, q\}^2$ ,  $\mathbf{m}_{|i}^{(k)}$  and  $\Sigma_{|i}^{(k)}$  are the mean vector and covariance matrix of the Gaussian vector ( $\mathbf{Z}_{-i}^{(k)} | Z_i^{(k)} = 0$ ), the index  $-i$  meaning that the  $i^{\text{th}}$  component is removed.

## 3 Gradient of the multipoint Expected Improvement

In this section, we provide an analytical formula for the gradient of  $q$ -EI. Getting such formula requires to carefully analyze the dependence of  $q$ -EI written in Eq. (4) on the batch locations  $\mathbf{X} \in \mathbb{R}^{q \times d}$ . This dependence is summarized in Fig. 1 and exhibits many chaining relations. In the forthcoming multivariate calculations, we use the following notations. Given two Banach spaces  $E$  and  $F$ , and a differentiable function  $g : E \rightarrow F$ , the differential of  $g$  at point  $x$ , written  $d_x[g] : E \rightarrow F$ , is the bounded linear map that best approximate  $g$  in the neighborhood of  $x$ . In the case where  $E = \mathbb{R}^p$  and  $F = \mathbb{R}$ , it is well known that  $\forall h \in E, d_x[g](h) = \langle \nabla g(x), h \rangle$ . More generally the differential can be written in terms of Jacobian matrices, matrix derivatives and/or matrix scalar products where  $E$  and/or  $F$  are  $\mathbb{R}^p$  or  $\mathbb{R}^{p \times p}$ . To simplify notations and handle the different indices in Eq. (4), we fix the indices  $i$  and  $k$  and focus on differentiating the function  $\text{EI}^{(k)(i)}$ , standing for the generic term of the double

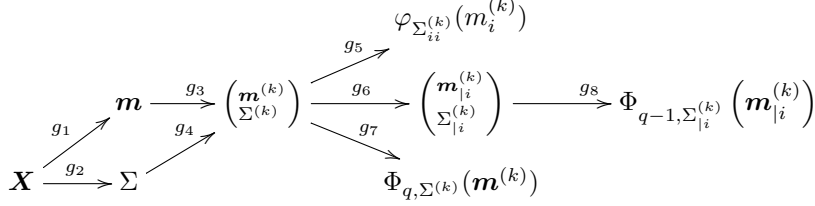


Figure 1: Link between the different terms of Eq. (4) and the batch of points  $\mathbf{X}$

sums in Eq. (4). We can perform the calculation of  $d_{\mathbf{X}} [\text{EI}^{(k)(i)}]$  by noticing than  $\text{EI}^{(k)(i)}$  can be rewritten using the functions  $g_j$ ,  $1 \leq j \leq 8$  defined on Fig. 1 as follows:

$$\text{EI}^{(k)(i)} = (m_k - T) \cdot g_7 \circ G + g_4 \circ g_2 \cdot g_5 \circ G \cdot g_8 \circ g_6 \circ G, \quad (5)$$

where  $G = (g_3 \circ g_1, g_4 \circ g_2)$ ,  $\circ$  is the composition operator and  $\cdot$  the multiplication operator. The differentiation then consists in applying classical differentiation formulas for products and compositions to Eq. (5). Proposition 3 summarizes the results. For conciseness, the formulae of the differentials involved in Eq. (6) are justified in the Appendix. The calculations notably rely on the differential of a normal cumulative distribution functions with respect to its covariance matrix obtained via Plackett's formula [4].

**Proposition 1.** *The differential of the multipoint Expected Improvement criterion of Eq. (4) is given by  $d_{\mathbf{X}} [\text{EI}] = \sum_{k=1}^q \sum_{i=1}^q d_{\mathbf{X}} [\text{EI}^{(k)(i)}]$ , with*

$$\begin{aligned} d_{\mathbf{X}} [\text{EI}^{(k)(i)}] &= d_{\mathbf{X}} [m_k] \cdot g_7 \circ G + (m_k - T) \cdot d_{G(\mathbf{X})} [g_7] \circ d_{\mathbf{X}} [G] \quad (6) \\ &\quad + d_{g_2(\mathbf{X})} [g_4] \circ d_{\mathbf{X}} [g_2] \cdot g_5 \circ G \cdot g_8 \circ g_6 \circ G \\ &\quad + g_4 \circ g_2 \cdot d_{G(\mathbf{X})} [g_5] \circ d_{\mathbf{X}} [G] \cdot g_8 \circ g_6 \circ G \\ &\quad + g_4 \circ g_2 \cdot g_5 \circ G \cdot d_{g_6(G(\mathbf{X}))} [g_8] \circ d_{G(\mathbf{X})} [g_6] \circ d_{\mathbf{X}} [G], \end{aligned}$$

where the  $g_j$ 's are the functions introduced in Fig. 1. The  $g_j$ 's and their respective differentials are as follow :

- $g_1 : \mathbf{X} \in D^q \rightarrow g_1(\mathbf{X}) = (\mu_n(\mathbf{x}_j))_{1 \leq j \leq q} \in \mathbb{R}^q$ ,  
 $d_{\mathbf{X}} [g_1](H) = (\langle \nabla \mu_n(\mathbf{x}_j), H_{j,1:d}^T \rangle)_{1 \leq j \leq q}$ ,  
with  $\nabla \mu_n(\mathbf{x}_j) = \nabla \mu(\mathbf{x}_j) + \left( \frac{\partial \mathbf{c}_n(\mathbf{x}_j)^T}{\partial \mathbf{x}_\ell} \right)_{1 \leq \ell \leq d} \mathbf{C}_n^{-1} (\mathbf{y}_{1:n} - \mu(\mathbf{x}_{1:n}))$ .
- $g_2 : \mathbf{X} \in D^q \rightarrow g_2(\mathbf{X}) = (C_n(\mathbf{x}_j, \mathbf{x}_\ell))_{1 \leq j, \ell \leq q} \in \mathcal{S}_{++}^q$ .  $\mathcal{S}_{++}^q$  is the set of  $q \times q$  positive definite matrices.

$$d_{\mathbf{X}} [g_2] (H) = \left( \left\langle \nabla_{\mathbf{x}} C_n(\mathbf{x}_j, \mathbf{x}_\ell), H_{j,1:d}^\top \right\rangle + \left\langle \nabla_{\mathbf{x}} C_n(\mathbf{x}_\ell, \mathbf{x}_j), H_{\ell,1:d}^\top \right\rangle \right)_{1 \leq j, \ell \leq q},$$

with  $\nabla_{\mathbf{x}} C_n(\mathbf{x}, \mathbf{x}') = \nabla_{\mathbf{x}} C(\mathbf{x}, \mathbf{x}') - \left( \frac{\partial \mathbf{c}_n(\mathbf{x})^\top}{\partial \mathbf{x}_p} \right)_{1 \leq p \leq d} \mathbf{C}_n^{-1} \mathbf{c}_n(\mathbf{x}')$ .

- $G : \mathbf{X} \rightarrow (\mathbf{m}^{(k)}, \Sigma^{(k)})$ ,  $d_{\mathbf{X}} [G] = (L^{(k)} d_{\mathbf{X}} [g_1], L^{(k)} d_{\mathbf{X}} [g_2] L^{(k)\top})$ .
- $g_7 : (\mathbf{a}, \Gamma) \in \mathbb{R}^q \times \mathcal{S}_{++}^q \rightarrow \Phi_{q,\Gamma}(\mathbf{a}) \in \mathbb{R}$ ,  
 $d_{G(\mathbf{X})} [g_7] (\mathbf{h}, H) = \langle \mathbf{h}, \nabla_{\mathbf{x}} \Phi_{q,\Sigma^{(k)}}(\mathbf{m}^{(k)}) \rangle + \text{tr}(H \nabla_{\Sigma} \Phi_{q,\Sigma^{(k)}}(\mathbf{m}^{(k)}))$ .  $\nabla_{\mathbf{x}} \Phi_{q,\Sigma^{(k)}}$   
and  $\nabla_{\Sigma} \Phi_{q,\Sigma^{(k)}}$  are the gradient of the multivariate Gaussian CDF with respect to  $\mathbf{x}$  and to the covariance matrix, given in appendix.
- $g_4 : \Sigma \rightarrow \Sigma^{(k)}$ ,  $d_{g_2(\mathbf{X})} [g_4] (H) = L^{(k)} H L^{(k)\top}$ .
- $g_5 : (\mathbf{a}, \Gamma) \in \mathbb{R}^q \times \mathcal{S}_{++}^q \rightarrow \varphi_{\Gamma_{ii}}(a_i) \in \mathbb{R}$ ,  
 $d_{G(\mathbf{X})} [g_5] (\mathbf{h}, H) = \left( -\frac{a_i}{\Gamma_{ii}} h_i + \frac{1}{2} \left( \frac{a_i^2}{\Gamma_{ii}^2} - \frac{1}{\Gamma_{ii}} \right) H_{ii} \right) \varphi_{\Gamma_{ii}}(a_i)$
- $g_6 : (\mathbf{m}^{(k)}, \Sigma^{(k)}) \in \mathbb{R}^q \times \mathcal{S}_{++}^q \rightarrow (\mathbf{m}_{|i}^{(k)}, \Sigma_{|i}^{(k)})$ ,

$$d_{(\mathbf{m}^{(k)}, \Sigma^{(k)})} [g_6] (h, H) = \left( \mathbf{h}_{-i} - \frac{\mathbf{h}_i}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)} + \frac{m_i^{(k)} H_{ii}}{\Sigma_{ii}^{(k)2}} \Sigma_{-i,i}^{(k)} - \frac{m_i^{(k)}}{\Sigma_{ii}^{(k)}} H_{-i,i}, \right.$$

$$\left. H_{-i,-i} + \frac{H_{ii}}{\Sigma_{ii}^{(k)2}} \Sigma_{-i,i}^{(k)} \Sigma_{-i,i}^{(k)\top} - \frac{1}{\Sigma_{ii}^{(k)}} H_{-i,i} \Sigma_{-i,i}^{(k)\top} - \frac{1}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)} H_{-i,i}^\top \right)$$

- $g_8 : (\mathbf{a}, \Gamma) \in \mathbb{R}^{q-1} \times \mathcal{S}_{++}^{q-1} \rightarrow \Phi_{q-1,\Gamma}(\mathbf{a}) \in \mathbb{R}$ ,  
 $d_{g_6(G(\mathbf{X}))} [g_8] = \langle \mathbf{h}, \nabla_{\mathbf{x}} \Phi_{q,\Sigma^{(k)}}(\mathbf{m}^{(k)}) \rangle + \text{tr}(H \nabla_{\Sigma} \Phi_{q,\Sigma^{(k)}}(\mathbf{m}^{(k)}))$ .

The gradient of  $q$ -EI, relying on Eq. (6) is implemented in the version 1.5 of the DiceOptim R package [9], together with a gradient-based local optimization algorithm. In the next section, we show that the analytical computation of the gradient offers substantial computational savings compared to numerical computation based on a finite-difference scheme. In addition, we investigate the performances of the batch-sequential EGO algorithm consisting in sequentially maximizing  $q$ -EI, and we compare it with the BUCB algorithm of [10].

## 4 Numerical tests

### 4.1 Computation time

In this section, we illustrate the benefits – in terms of computation time – of using the analytical gradient formula of Section 3. We compare computation times of gradients computed analytically and numerically, through finite differences schemes. It is important to note that the computation of both  $q$ -EI and

	$\Phi_{q-3}$	$\Phi_{q-2}$	$\Phi_{q-1}$	$\Phi_q$	Total
analytic $q$ -EI	0	0	$q^2$	$q$	$O(q^2)$
finite differences gradient	0	0	$q(d+1)q^2$	$q(d+1)q$	$O(dq^3)$
analytic gradient	$q^2 \frac{q-1}{2}$	$q \frac{q-1}{2} + q^3$	$q^2 + 2q^2$	$q$	$O(q^4)$

Table 1: Total number of calls to the CDF of the multivariate Gaussian distribution for computing  $q$ -EI or its gradient for a function with  $d$  input variables. The last column gives the computational complexity as a function of  $q$  and  $d$ .

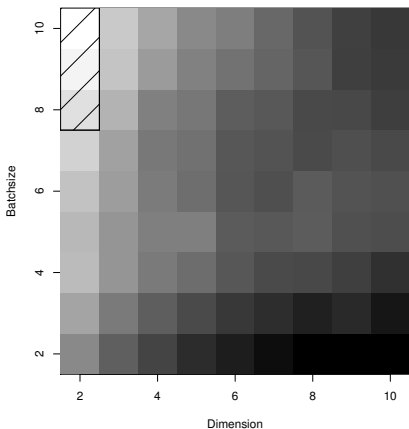


Figure 2: Ratio between computation times of the numerical and analytical gradient of  $q$ -EI as a function of the dimension  $d$  and the batch size  $q$ . The hatched area indicates a ratio below 1.

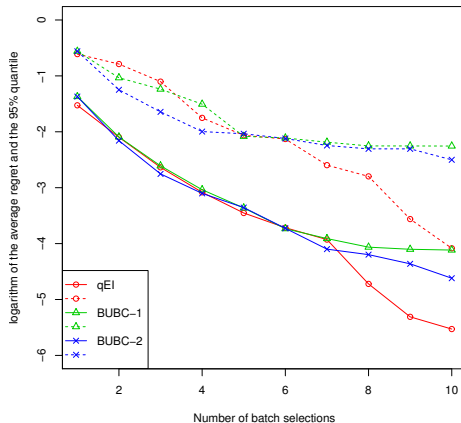


Figure 3: Logarithm of the average (plain lines) and 95% quantile (dotted lines) of the regret for three different batch-sequential optimization strategies (see Section 4.2 for detail).

its gradient (see, Eqs. (4),(6)) involve several calls to the cumulative distribution functions (CDF) of the multivariate normal distribution. The latter CDF is computed numerically with the algorithms of [13] wrapped in the `mnormt` R package [2]. In our implementation, computing this CDF turns out to be the main bottleneck in terms of computation time. The total number of calls to this CDF (be it in dimension  $q$ ,  $q-1$ ,  $q-2$  or  $q-3$ ) is summarized in Table 1. From this table, let us remark that the number of CDF calls does not depend on  $d$  for the analytical  $q$ -EI gradient and is proportional to  $d$  for the numerical gradient. The use of the analytical gradient is thus expected to bring savings when  $q$  is not too large compared to  $d$ . Figure 2 depicts the ratio of computation times between numerical and analytical gradient, as a function of  $q$  and  $d$ . These were obtained by averaging the evaluation times of  $q$ -EI’s gradient at 10 randomly-generated batches of size  $q$  for a given function in dimension  $d$  being a sample path of a GP with separable Matérn(3/2) covariance function [21]. In the next section, we use the values  $q = 6$  and  $d = 5$  and we rely exclusively on the analytical  $q$ -EI formula which we now know to be faster.



## 4.2 Tests

### 4.2.1 Experimental setup

We now compare the performances of two parallel Bayesian optimization algorithm based, respectively, on the UCB approach of [23] and on sequential  $q$ -EI maximizations. We consider a minimization problem in dimension  $d = 5$  where  $n = 50$  evaluations are performed initially and 10 batches of  $q = 6$  observations are sequentially added. The objective functions are 50 different sample realizations of a zero mean GP with unit variance and separable isotropic Matérn(3/2) covariance function with range parameter equal to one. Both algorithms use the same initial design of experiment of  $n$  points which are all S-optimal random Latin Hypercube designs [17]. The mean and covariance function of the underlying GP are supposed to be known. Since it is difficult to draw sample realizations of the GP on the whole input space  $D := [0, 1]^d$ , we instead draw 50 samples on a set of 2000 space-filling locations and interpolate each sample in order to obtain the 50 objective functions.

Two variants of the BUCB algorithms are tested. Each of them constructs a batch by sequentially minimizing the kriging quantile  $\mu_n^*(\mathbf{x}) - \beta_n s_n(\mathbf{x})$  where  $s_n(\mathbf{x}) = \sqrt{C_n(\mathbf{x}, \mathbf{x})}$  is the posterior standard deviation at step  $n$  and  $\mu_n^*(\mathbf{x})$  is the posterior mean conditioned both on the response at previous points and at points already selected in the current batch, with a dummy response fixed to their posterior means in the latter case. Following the settings of [10], in the first and second variant of BUCB, the coefficients  $\beta_n$  are given by:

$$\beta_n^{(1)} := 2\beta_{\text{mult}} \log \left( \frac{\pi^2 d}{6\delta} (k+1)^2 \right) \quad \text{and} \quad \beta_n^{(2)} := 2\beta_{\text{mult}} \log \left( \frac{\pi^2 d}{6\delta} (1+qk)^2 \right) \quad (7)$$

where  $\beta_{\text{mult}} = 0.1$ ,  $\delta = 0.1$ , and  $k$  is the number of already evaluated batches at time  $n$ , i.e., here,  $k \in \{0, \dots, 9\}$ . The BUCB1 strategy is expected to select locations in regions with low posterior mean (exploitation) while BUCB2 is meant to favour more exploration due to a larger  $\beta_n$ . The minimization of the kriging quantile presented above is performed using a genetic algorithm [18]. Regarding the algorithm based on  $q$ -EI sequential maximization, we propose to use a multi-start BFGS algorithm with analytical gradient. This algorithm operates gradient descents directly in the space of dimension  $qd = 30$ . To limit computation time, the number of starting batches in the multi-start is set to 3. These 3 batches are obtained by running the BUCB1 algorithm presented above with 3 different values of  $\beta_{\text{mult}}$  equals to 0.05, 0.1, 0.2 respectively.

At each iteration, we measure the regrets of each algorithm and average them over the 50 experiments. To facilitate the interpretation of results, we first focus on the results of the algorithms after 1 iteration, i.e. after having added only 1 batch of  $q$  points. We then discuss the results when 10 iterations are run.

### 4.2.2 First step of the optimization

To start with, we focus on the selection of the first batch. Table 2 compares the average  $q$ -EI and real improvement obtained for the three selection rules. For the first iteration only, the BUCB1 and BUCB2 selection rules are exactly the same. Since  $q$ -EI is the one-step optimal, it is not a surprise that it performs

Selection rule	Average expected improvement ( $q$ -EI)	Average realized improvement
$q$ -EI	0.672	0.697
BUCB	0.638	0.638

Table 2: Expected and observed first batch Improvement for  $q$ -EI and BUCB batch selection methods, in average for 50 functions.

better at iteration 1 with our settings where the objective functions are sample realizations of a GP. If only one iteration is performed, improving the  $q$ -EI is equivalent to improving the average performance. However, we point out that, in application, the maximization of  $q$ -EI was not straightforward. It turns out that the batches proposed by the BUCB algorithms were excellent initial candidates in our descent algorithms. The use of other rules for the starting batches, with points sampled uniformly or according to a density proportional to the one-point EI, did not manage to yield this level of performance.

### 4.2.3 10 optimization steps

The average regret of the different batch selection rules over 10 iteration is depicted in Fig. 3. This Figure illustrates that choosing the one-step optimal criterion is not necessarily optimal if more than one iteration is run [14]. Indeed, after two steps,  $q$ -EI maximization is already beaten by BUCB2, and  $q$ -EI becomes better again after iteration 7. Among the 50 optimized functions,  $q$ -EI maximization gives the smallest 10-steps final regret for only 30% of functions, against 52% for the BUCB1 and 18% for the BUCB2. On the other hand, the  $q$ -EI selection rule is eventually better in average since, for some functions, BUCB is beaten by  $q$ -EI by a wide margin. This is further illustrated with the curve of the 95% quantile of the regret which indicates that, for the worst simulations,  $q$ -EI performs better. This gain in robustness alone explains the better average performance of  $q$ -EI. Such improved performance comes at a price : the computational time of our multistart BFGS algorithm with analytical gradient is 4.1 times higher compared to the BUCB computation times.

## 5 Conclusion

In this article, we give a closed-form expression of the gradient of the multi-point Expected Improvement criterion, enabling an efficient  $q$ -EI maximization

at reduced computational cost. Parallel optimization strategies based on maximization of  $q$ -EI have been tested and are ready to be used on real test case with the DiceOptim R package. The BUCB algorithm turns out to be a good competitor to  $q$ -EI maximization, with a lower computational cost, and also gives good starting batches for the proposed multistart BFGS algorithm. In general, however, the maximization of  $q$ -EI remains a difficult problem. An interesting perspective is to develop algorithms taking advantage of some particular properties of the  $q$ -EI function in the space of dimension  $qd$ , for example its invariance to point permutations. Other research perspectives include deriving cheap but trustworthy approximations of  $q$ -EI and its gradient. Finally, as illustrated in the application,  $q$ -EI sequential maximizations have no reason to constitute optimal decisions for a horizon beyond one batch. Although the optimal policy is known [14], its implementation in practice remains an open problem.

### Acknowledgement

Part of this work has been conducted within the frame of the ReDice Consortium, gathering industrial (CEA, EDF, IFPEN, IRSN, Renault) and academic (École des Mines de Saint-Étienne, INRIA, and the University of Bern) partners around advanced methods for Computer Experiments.

## 6 Appendix : Differential calculus

- $g_1$  and  $g_2$  are functions giving respectively the mean of  $\mathbf{Y}(\mathbf{X})$  and its covariance. Each component of these functions is either a linear or a quadratic combination of the trend function  $\boldsymbol{\mu}$  or the covariance function  $C$  evaluated at different points of  $\mathbf{X}$ . The results are obtained by matrix differentiation. See the appendix B of [21] for a similar calculus.
- $g_3$  (resp.  $g_4$ ) is the affine (resp. linear) transformation of the mean vector  $\mathbf{m}$  into  $\mathbf{m}^{(k)}$  (resp. the covariance matrix  $\Sigma$  into  $\Sigma^{(k)}$ ). The differentials are then expressed in terms of the same linear transformation :

$$d_{\mathbf{m}} [g_3] (\mathbf{h}) = L^{(k)} \mathbf{h} \quad \text{and} \quad d_{\Sigma} [g_4] (H) = L^{(k)} H L^{(k)\top}.$$

- $g_5$  is defined by  $g_5(\mathbf{m}^{(k)}, \Sigma^{(k)}) = \varphi_{\Sigma_{ii}^{(k)}}(m_i^{(k)})$ . Then the result is obtained by differentiating the univariate Gaussian probability density function with respect to its mean and variance parameters. Indeed we have :

$$d_{(\mathbf{m}^{(k)}, \Sigma^{(k)})} [g_5] (h, H) = d_{\mathbf{m}^{(k)}} [g_5(\cdot, \Sigma^{(k)})] (h) + d_{\Sigma^{(k)}} [g_5(\mathbf{m}^{(k)}, \cdot)] (H)$$

- $g_6$  gives the mean and the covariance of  $\mathbf{Z}_{-i}^{(k)} | Z_i = 0$ . We have :

$$\left( \mathbf{m}_{-i}^{(k)}, \Sigma_{-i}^{(k)} \right) = g_6 \left( \mathbf{m}^{(k)}, \Sigma^{(k)} \right) = \left( \mathbf{m}_{-i}^{(k)} - \frac{m_i^{(k)}}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)}, \Sigma_{-i,-i}^{(k)} - \frac{1}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)} \Sigma_{-i,i}^{(k)\top} \right)$$

$$d_{(\mathbf{m}^{(k)}, \Sigma^{(k)})} [g_6] (\mathbf{h}, H) = d_{\mathbf{m}^{(k)}} \left[ g_6 \left( \cdot, \Sigma^{(k)} \right) \right] (\mathbf{h}) + d_{\Sigma^{(k)}} [g_6] \left( \mathbf{m}^{(k)}, \cdot \right) (H),$$

$$\text{with : } d_{\mathbf{m}^{(k)}} \left[ g_6 \left( \cdot, \Sigma^{(k)} \right) \right] (\mathbf{h}) = \left( \mathbf{h}_{-i} - \frac{\mathbf{h}_i}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)}, 0 \right)$$

$$\text{and : } d_{\Sigma^{(k)}} \left[ g_6 \left( \mathbf{m}^{(k)}, \cdot \right) \right] (H) = \left( \frac{m_i^{(k)} H_{ii}}{\Sigma_{ii}^{(k)2}} \Sigma_{-i,i}^{(k)} - \frac{m_i^{(k)}}{\Sigma_{ii}^{(k)}} H_{-i,i}, \right.$$

$$\left. H_{-i,-i} + \frac{H_{ii}}{\Sigma_{ii}^{(k)2}} \Sigma_{-i,i}^{(k)} \Sigma_{-i,i}^{(k)\top} - \frac{1}{\Sigma_{ii}^{(k)}} H_{-i,i} \Sigma_{-i,i}^{(k)\top} - \frac{1}{\Sigma_{ii}^{(k)}} \Sigma_{-i,i}^{(k)} H_{-i,i}^\top \right)$$

- $g_7$  and  $g_8$  : these two functions take a mean vector and a covariance matrix in argument and give a probability in output :  $\Phi_{q, \Sigma^{(k)}} (-\mathbf{m}^{(k)}) = g_7 (\mathbf{m}^{(k)}, \Sigma^{(k)})$ ,  $\Phi_{q-1, \Sigma_{|i}^{(k)}} (-\mathbf{m}_{|i}^{(k)}) = g_8 (\mathbf{m}_{|i}^{(k)}, \Sigma_{|i}^{(k)})$ . So, for  $\{p, \Gamma, \mathbf{a}\} = \{q, \Sigma^{(k)}, -\mathbf{m}^{(k)}\}$  or  $\{q-1, \Sigma_{|i}^{(k)}, -\mathbf{m}_{|i}^{(k)}\}$ , we face the problem of differentiating a function  $\Phi : (\mathbf{a}, \Gamma) \rightarrow \Phi_{p, \Gamma}(\mathbf{a})$ , with respect to  $(\mathbf{a}, \Gamma) \in \mathbb{R}^p \times \mathcal{S}_{++}^p$ :

$$d_{(\mathbf{a}, \Gamma)} [\Phi] (\mathbf{h}, H) = d_{\mathbf{a}} [\Phi(\cdot, \Gamma)] (\mathbf{h}) + d_{\Gamma} [\Phi(\mathbf{a}, \cdot)] (H).$$

The the first differential of this sum can be written :

$$d_{\mathbf{a}} [\Phi(\cdot, \Gamma)] (\mathbf{h}) = \left\langle \left( \frac{\partial}{\partial a_i} \Phi(\mathbf{a}, \Gamma) \right)_{1 \leq i \leq p}, \mathbf{h} \right\rangle,$$

$$\text{with : } \frac{\partial}{\partial a_i} \Phi(\mathbf{a}, \Gamma) = \int_{-\infty}^{a_1} \dots \int_{-\infty}^{a_{i-1}} \int_{-\infty}^{a_{i+1}} \dots \int_{-\infty}^{a_p} \varphi_{p, \Gamma}(u_{-i}, a_i) \mathbf{d}\mathbf{u}_{-i} = \varphi_{1, \Gamma_{ii}} \Phi_{p-1, \Gamma_{|i}}(\mathbf{a}_{|i}).$$

The last equality is obtained with the identity :  $\forall \mathbf{u} \in \mathbb{R}^q$ ,  $\varphi_{q, \Gamma}(\mathbf{u}) = \varphi_{1, \Gamma_{ii}}(u_i) \varphi_{p-1, \Gamma_{|i}}(\mathbf{u}_{|i})$ , with  $\mathbf{u}_{|i} = \mathbf{u}_{-i} - \frac{u_i}{\Gamma_{ii}} \boldsymbol{\Gamma}_{-i,i}$  and  $\Gamma_{|i} = \Gamma_{-i,-i} - \frac{1}{\Gamma_{ii}} \boldsymbol{\Gamma}_{-i,i} \boldsymbol{\Gamma}_{-i,i}^\top$ . The second differential is :

$$d_{\Gamma} [\Phi(\mathbf{a}, \cdot)] (H) := \frac{1}{2} \text{tr} \left( H \cdot \left( \frac{\partial \Phi}{\partial \Gamma_{ij}}(\mathbf{a}, \Gamma) \right)_{i,j \leq p} \right) = \frac{1}{2} \text{tr} \left( H \cdot \left( \frac{\partial^2 \Phi}{\partial a_i \partial a_j}(\mathbf{a}, \Gamma) \right)_{i,j \leq p} \right)$$

$$\text{where : } \frac{\partial^2 \Phi}{\partial a_i \partial a_j}(\mathbf{a}, \Gamma) = \begin{cases} \varphi_{2, \Sigma_{\{i,j\}, \{i,j\}}}(\mathbf{x}_i, \mathbf{x}_j) \Phi_{p-2, \Sigma_{|ij}}(\mathbf{x}_{|ij}), & \text{if } i \neq j, \\ -\frac{x_i}{\Gamma_{ii}} \frac{\partial}{\partial a_i} \Phi_{\Gamma}(\mathbf{a}, \Gamma) - \sum_{\substack{j=1 \\ j \neq i}}^p \frac{1}{\Gamma_{ii}} \Gamma_{ij} \frac{\partial^2}{\partial a_i \partial a_j} \Phi(\mathbf{a}, \Gamma). \end{cases}$$

## References

- [1] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [2] A. Azzalini and A. Genz. *The R package mnormt: The multivariate normal and t distributions (version 1.5-1)*, 2014.

- [3] J. Bect, D. Ginsbourger, L. Li, V. Picheny, and E. Vazquez. Sequential design of computer experiments for the estimation of a probability of failure. *Statistics and Computing*, 22(3):773–793, 2011.
- [4] S. M. Berman. An extension of placketts differential equation for the multivariate normal density. *SIAM Journal on Algebraic Discrete Methods*, 8(2):196–197, 1987.
- [5] E. Brochu, 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. eprint arXiv:1012.2599, arXiv.org, December 2010.
- [6] A. Bull. Convergence rates of efficient global optimization algorithms. *Journal of Machine Learning Research*, 12:2879–2904, 2011.
- [7] C. Chevalier. *Fast uncertainty reduction strategies relying on Gaussian process models*. PhD thesis, University of Bern, 2013.
- [8] C. Chevalier and D. Ginsbourger. *Learning and Intelligent Optimization - 7th International Conference, Lion 7, Catania, Italy, January 7-11, 2013, Revised Selected Papers*, chapter fast computation of the multipoint expected improvement with applications in batch selection, pages 59-69. Springer, 2014.
- [9] D. Ginsbourger and V. Picheny and O. Roustant and with contributions by C. Chevalier and S. Marmin and T. Wagner. *DiceOptim: Kriging-Based Optimization for Computer Experiments*, 2015. R package version 1.5.
- [10] T. Desautels, A. Krause, and J. Burdick. Parallelizing exploration-exploitation tradeoffs with gaussian process bandit optimization. In *ICML*, 2012.
- [11] P. I. Frazier. Parallel global optimization using an improved multi-points expected improvement criterion. In *INFORMS Optimization Society Conference, Miami FL*, 2012.
- [12] P. I. Frazier, W. B. Powell, and S. Dayanik. A knowledge-gradient policy for sequential information collection. *SIAM Journal on Control and Optimization*, 47(5):2410–2439, 2008.
- [13] A. Genz. Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–149, 1992.
- [14] D. Ginsbourger and R. Le Riche. Towards gaussian process-based optimization with finite time horizon. In Alessandra Giovagnoli, Anthony C. Atkinson, Bernard Torsney, and Caterina May, editors, *mODa 9 96 Advances in Model-Oriented Design and Analysis*, Contributions to Statistics, pages 89–96. Physica-Verlag HD, 2010.

- [15] D. Ginsbourger, R. Le Riche, and L. Carraro. Kriging is well-suited to parallelize optimization. In *Computational Intelligence in Expensive Optimization Problems*, volume 2 of *Adaptation Learning and Optimization*, pages 131–162. Springer, 2010.
- [16] D. R. Jones, M. Schonlau, and J. William. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998.
- [17] Q. Y. Kenny, W. Li, and A. Sudjianto. Algorithmic construction of optimal symmetric latin hypercube designs. *Journal of statistical planning and inference*, 90(1):145–159, 2000.
- [18] W. Mebane and J. Sekhon. Genetic optimization using derivatives: The rgenoud package for r. *Journal of Statistical Software*, Vol. 42, Issue 11:1–26, 2011.
- [19] J. Mockus, V. Tiesis, and A. Zilinskas. The application of Bayesian methods for seeking the extremum. In L. Dixon and Eds G. Szego, editors, *Towards Global Optimization*, volume 2, pages 117–129. Elsevier, 1978.
- [20] C. R. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [21] O. Roustant, D. Ginsbourger, and Y. Deville. DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by Kriging-Based Metamodelling and Optimization. *Journal of Statistical Software*, 51 (1):1–55, 2012.
- [22] M. Schonlau. *Computer Experiments and global optimization*. PhD thesis, University of Waterloo, 1997.
- [23] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Information-theoretic regret bounds for gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 58(5):3250–3265, 2012.
- [24] E. Vazquez and J. Bect. Convergence properties of the expected improvement algorithm with fixed mean and covariance functions. *Journal of Statistical Planning and inference*, 140(11):3088–3095, 2010.
- [25] J. Villemonteix, E. Vazquez, and E. Walter. An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509–534, 2009.