

# OPTIMIZING POSTERIOR SAMPLES FOR BAYESIAN OPTIMIZATION VIA ROOTFINDING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Bayesian optimization devolves the global optimization of a costly objective function to the global optimization of a sequence of acquisition functions. This inner-loop optimization can be catastrophically difficult if it involves posterior sample paths, especially in higher dimensions. We introduce an efficient global optimization strategy for posterior samples based on global rootfinding. It provides gradient-based optimizers with [two sets of judiciously selected starting points](#), designed to combine exploration and exploitation. [The number of starting points can be kept small without sacrificing optimization quality](#). Remarkably, [even with just one point from each set, the global optimum is discovered most of the time](#). The algorithm scales practically linearly to high dimensions, [breaking the curse of dimensionality](#). For Gaussian process Thompson sampling (GP-TS), we demonstrate remarkable improvement in both inner- and outer-loop optimization, [surprisingly outperforming alternatives like EI and GP-UCB in most cases](#). Our approach also improves the performance of other posterior sample-based acquisition functions, [such as variants of entropy search](#). Furthermore, we propose a sample-average formulation of GP-TS, which has a parameter to explicitly control exploitation and can be computed at the cost of one posterior sample.

## 1 INTRODUCTION

Bayesian optimization (BO) is a highly successful approach to the global optimization of expensive-to-evaluate black-box functions, with applications ranging from hyper-parameter training of machine learning models to scientific discovery and engineering design (Jones et al., 1998; Snoek et al., 2012; Frazier, 2018; Garnett, 2023). Many BO strategies are also backed by strong theoretical guarantees on their convergence to the global optimum (Srinivas et al., 2010; Bull, 2011; Russo & Van Roy, 2014; Chowdhury & Gopalan, 2017).

Consider the global optimization problem  $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  where  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$  represents the vector of input variables and  $f(\mathbf{x}) \in \mathbb{R}$  the objective function which can be evaluated at a significant cost, subject to observation noise. At its core, BO is a sequential optimization algorithm that uses a probabilistic model of the objective function to guide its evaluation decisions. Starting with a prior probabilistic model and some initial data, BO derives an acquisition function  $\alpha(\mathbf{x})$  from the posterior model, which is much easier to evaluate than the objective function and often has easy-to-evaluate derivatives. The acquisition function is then optimized globally, using off-the-shelf optimizers, to provide a location to evaluate the objective function. This process is iterated until some predefined stopping criteria are met.

Effectively there are two nested iterations in BO: the outer-loop seeks to optimize the objective function  $f(\mathbf{x})$ , and the inner-loop seeks to optimize the acquisition function  $\alpha(\mathbf{x})$  at each BO iteration. The premise of BO is that the inner-loop optimization can be solved accurately and efficiently, so that the outer-loop optimization proceeds informatively with a negligible added cost. In fact, the convergence guarantees of many BO strategies assume *exact* global optimization of the acquisition function. However, the efficient and accurate global optimization of acquisition functions is less trivial than it is often assumed to be (Wilson et al., 2018).

Acquisition functions are, in general, highly non-convex and have many local optima. In addition, many common acquisition functions are mostly flat surfaces with a few peaks (Rana et al., 2017),

054 which take up an overwhelmingly large portion of the domain as the input dimension grows. This  
 055 creates a significant challenge for generic global optimization methods.

056 Some acquisition functions involve sample functions from the posterior model, which need to be  
 057 optimized globally. Gaussian process Thompson sampling (GP-TS) (Chowdhury & Gopalan, 2017)  
 058 uses posterior sample paths directly as random acquisition functions. In many information-theoretic  
 059 acquisition functions such as [entropy search \(ES\)](#) (Hennig & Schuler, 2012), predictive entropy  
 060 search (PES) (Hernández-Lobato et al., 2014), max-value entropy search (MES) (Wang & Jegelka,  
 061 2017), and joint entropy search (JES) (Tu et al., 2022; Hvarfner et al., 2022), multiple posterior  
 062 samples are drawn and optimized to find their global optimum location and/or value. Such acqui-  
 063 sition functions are celebrated for their nice properties in BO: TS has strong theoretical guarantees  
 064 (Russo & Van Roy, 2014; 2016) and can be scaled to high dimensions (Mutny & Krause, 2018);  
 065 information-theoretic acquisition functions are grounded in principles for optimal experimental design  
 066 (MacKay, 2003); and both types can be easily parallelized in synchronous batches (Shah &  
 067 Ghahramani, 2015; Hernández-Lobato et al., 2017; Kandasamy et al., 2018). However, posterior  
 068 sample paths are much more complex than other acquisition functions, as they fluctuate throughout  
 069 the design space, and are less smooth than the posterior mean and marginal variance. The latter  
 070 are the basis of many acquisition functions, such as expected improvement (EI) (Jones et al., 1998),  
 071 probability of improvement (PI) (Kushner, 1964), and upper confidence bound (GP-UCB) (Srini-  
 072 vas et al., 2010). As a consequence, posterior sample paths have many more local optima, and the  
 073 number scales exponentially with the input dimension.

074 While there is a rich literature on prior probabilistic models and acquisition functions for BO, global  
 075 optimization algorithms for acquisition functions have received little attention. One class of global  
 076 optimization methods is derivative-free, such as the dividing rectangles (DIRECT) algorithm (Jones  
 077 et al., 1993), covariance matrix adaptation evolution strategy (CMA-ES) algorithms (Hansen et al.,  
 078 2003), and genetic algorithms (Mitchell, 1998). Gradient-based multistart optimization, on the other  
 079 hand, is often seen as the best practice to reduce the risk of getting trapped in local minima (Kim &  
 080 Choi, 2021), and enjoys the efficiency of being embarrassingly parallelizable. For posterior samples,  
 081 their global optimization may use joint sampling on a finite set of points (Kandasamy et al., 2018), or  
 082 approximate sampling of function realizations followed by gradient-based optimization (Hernández-  
 083 Lobato et al., 2014; Mutny & Krause, 2018). The selection of starting points is crucial for the success  
 084 of gradient-based multistart optimization. This selection can be deterministic (e.g., grid search),  
 random (Bergstra & Bengio, 2012; Balandat et al., 2020), or adaptive (Feo & Resende, 1995).

085 In this paper, we propose an adaptive strategy for selecting starting points for gradient-based mul-  
 086 tistart optimization of posterior samples. This algorithm builds on the decomposition of poste-  
 087 rior samples by pathwise conditioning, taps into robust software in univariate function computation  
 088 based on univariate global rootfinding, and exploits the separability of multivariate GP priors. Our  
 089 key contributions include:

- 090 • A novel strategy for the global optimization of posterior sample paths. The starting points are  
 091 dependent on the posterior sample, so that each is close to a local optimum that is a candidate for  
 092 the global optimum. The selection algorithm scales linearly to high dimensions.
- 093 • We give empirical results across a diverse set of problems with input dimensions ranging from 2  
 094 to 16, establishing the effectiveness of our optimization strategy. Although our algorithm is pro-  
 095 posed for the inner-loop optimization of posterior samples, perhaps surprisingly, we see significant  
 096 improvement in outer-loop optimization performance, which often allows acquisition functions  
 097 based on posterior samples to converge faster than other common acquisition functions.
- 098 • A new acquisition function via the posterior sample average that explicitly controls the  
 099 exploration–exploitation balance (Chapelle & Li, 2011), and can be generated at the same cost  
 100 as one posterior sample.

## 101 2 GENERAL BACKGROUND

102 **Gaussian Processes.** Consider an unknown  $d$ -dimensional function  $f_{\text{true}} : \mathcal{X} \mapsto \mathbb{R}$ , where domain  $\mathcal{X} \subset \mathbb{R}^d$ . We can  
 103 collect noisy observations of the function through the model  $y^i = f_{\text{true}}(\mathbf{x}^i) + \varepsilon^i$ ,  $i \in \{1, \dots, n\}$ ,  
 104 with  $\varepsilon \sim \mathcal{N}_n(0, \Sigma)$ . To model the function  $f_{\text{true}}$ , we use a Gaussian process (GP) as the prior  
 105 probabilistic model:  $f \sim \pi \in \mathcal{GP}$ . A GP is a random function  $f$  such that for any finite set of  
 106 points  $X = \{\mathbf{x}^i\}_{i=1}^n$ ,  $n \in \mathbb{N}$ , the values  $\mathbf{f}_n = (f(\mathbf{x}^i))_{i=1}^n$  have a multivariate Gaussian distribution  
 107

$\mathcal{N}_n(\boldsymbol{\mu}_n, \mathbf{K}_{n,n})$ , with mean  $\boldsymbol{\mu}_n = (\mu(\mathbf{x}^i))_{i=1}^n$  and covariance  $\mathbf{K}_{n,n} = [\kappa(\mathbf{x}^i, \mathbf{x}^j)]_{i,j \in n}^n$ . Here,  $\mu(\mathbf{x})$  is the mean function and  $\kappa(\mathbf{x}, \mathbf{x}')$  is the covariance function.

**Decoupled Representation of GP Posteriors.** Given a dataset  $\mathcal{D} = \{(\mathbf{x}^i, y^i)\}_{i=1}^n$ , the posterior model  $f|\mathcal{D}$  is also a GP. Samples from the posterior have a decoupled representation called pathwise conditioning, originally proposed in (Wilson et al., 2020; 2021):

$$(f|\mathcal{D})(\mathbf{x}) \stackrel{d}{=} f(\mathbf{x}) + \boldsymbol{\kappa}_{\cdot,n}(\mathbf{x})(\mathbf{K}_{n,n} + \boldsymbol{\Sigma})^{-1}(\mathbf{y} - \mathbf{f}_n - \boldsymbol{\varepsilon}), \quad f \sim \pi, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \boldsymbol{\Sigma}), \quad (1)$$

where  $f(\mathbf{x})$  is a sample path from the GP prior,  $\boldsymbol{\kappa}_{\cdot,n}(\mathbf{x}) = (\kappa(\mathbf{x}, \mathbf{x}^i))_{i=1}^n$  is the canonical basis,  $\mathbf{f}_n = (f(\mathbf{x}^i))_{i=1}^n$ , and  $\boldsymbol{\varepsilon}$  is a sample of the noise. We may interpret  $\mathbf{f}_n + \boldsymbol{\varepsilon}$  as a sample from the prior distribution of the observations  $\mathbf{y} = (y^i)_{i=1}^n$ . This representation has its roots in Matheron’s update rule that transforms a joint distribution of Gaussian variables into a conditional one (see e.g., Hoffman & Ribak (1991)). This formula is exact, in that  $\stackrel{d}{=}$  denotes equality in distribution, and it preserves the differentiability of the prior sample. It is also computationally efficient for posterior sampling: the cost is independent of the input dimension  $d$ , linear in the data size  $n$  at evaluation time, and the weight vector for  $\boldsymbol{\kappa}_{\cdot,n}(\mathbf{x})$  can be solved accurately using an iterative algorithm that scales linearly with  $n$  (Lin et al., 2023).

### 3 GLOBAL OPTIMIZATION OF POSTERIOR SAMPLE PATHS

In this section, we introduce an efficient algorithm for the global optimization of posterior sample paths. For this, we exploit the separability of prior samples and useful properties of posterior samples to judiciously select a set of starting points for gradient-based multistart optimizers.

**Assumptions.** Following Section 2, we impose a few common assumptions throughout this paper: (1) the domain is a hypercube:  $\mathcal{X} = \prod_{i=1}^d [x_i, \bar{x}_i]$ ; (2) prior covariance is separable:  $\kappa(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d \kappa_i(x_i, x'_i)$ ; (3) prior samples are continuously differentiable:  $f(\mathbf{x}; \omega) \in C^1$ . Without loss of generality, we also assume that the prior mean  $\mu(\mathbf{x}) = 0$ : any non-zero mean function can be subtracted from the data by replacing  $f_{\text{true}}$  with  $f_{\text{true}} - \mu$ . While additive and multiplicative compositions of univariate kernels can be used in the prior (Duvenaud et al., 2013), assumption (2) is the most popular choice in BO. It is possible to extend our method to generalized additive models.

#### 3.1 TS-ROOTS ALGORITHM

We observe that, given the assumptions, a posterior sample in eq. (1) can be rewritten as:

$$\tilde{f}(\mathbf{x}; \tilde{\omega}) = f(\mathbf{x}; \omega) + b(\mathbf{x}; \tilde{\omega}), \quad f(\mathbf{x}; \omega) = \prod_{i=1}^d f_i(x_i; \omega_i), \quad b(\mathbf{x}; \tilde{\omega}) = \sum_{j=1}^n v_j \kappa(\mathbf{x}, \mathbf{x}^j). \quad (2)$$

Here, the prior sample  $f(\mathbf{x})$  is a separable function determined by the random bits  $\omega$ . Data adjustment  $b(\mathbf{x})$  is a sum of the canonical basis with coefficients  $\mathbf{v} = (\mathbf{K}_{n,n} + \boldsymbol{\Sigma})^{-1}(\mathbf{y} - \mathbf{f}_n - \boldsymbol{\varepsilon})$ . Both the data adjustment and the posterior sample are determined by the random bits  $\tilde{\omega} = (\omega, \boldsymbol{\varepsilon})$ . In the following, we denote the prior and the posterior samples as  $f_\omega$  and  $\tilde{f}_{\tilde{\omega}}$ , respectively. **Our goal is to find the global minimum  $(\tilde{\mathbf{x}}_{\tilde{\omega}}^*, \tilde{f}_{\tilde{\omega}}^*)$  of the posterior sample  $\tilde{f}_{\tilde{\omega}}(\mathbf{x})$ .**

The global minimization of a generic function, in principle, requires finding all its local minima and then selecting the best among them. However, computationally efficient approaches to this problem are lacking even in low dimensions and, more pessimistically, the number of local minima grows exponentially as the domain dimension increases. The core idea of this work is to use the prior sample  $f_\omega$  as a surrogate of the posterior sample  $\tilde{f}_{\tilde{\omega}}$  for global optimization. Another key is to exploit the separability of the prior sample for efficient representation and ordering of its local extrema.

We define *TS-roots* as a global optimization algorithm for GP posterior samples (given the assumptions) via gradient-based multistart optimization. The starting points include: (1) a subset  $S_e$  of the local minima  $\tilde{X}_\omega$  of a corresponding prior sample; and (2) a subset  $S_x$  of the observed locations  $X$ . We call  $S_e$  the exploration set and  $S_x$  the exploitation set. Specifically, let

$$S_0 = \arg \operatorname{mink}_{\mathbf{x} \in \tilde{X}_\omega} (f_\omega(\mathbf{x}), n_0) \quad (3)$$

**Algorithm 1** TS-roots: Optimizing posterior samples via rootfinding

**Input:** prior samples  $f(\mathbf{x})$ ,  $\epsilon$ ; prior covariances  $\kappa(\mathbf{x}, \mathbf{x}')$ ,  $\Sigma$ ; dataset  $\mathcal{D}$ ; set sizes  $n_o, n_e, n_x$ .  
1:  $S_o \leftarrow \text{minsort}(f(\mathbf{x}), n_o)$   $\triangleright$  Smallest  $n_o$  local minima of the prior sample (Algorithm 4)  
2:  $[\tilde{\mathbf{x}}_e, I_e] \leftarrow \text{mink}(f(S_o), n_e)$   $\triangleright$  Smallest  $n_e$  values and indices of the posterior sample in  $S_o$   
3:  $[\tilde{\mathbf{x}}_x, I_x] \leftarrow \text{mink}(f(X), n_x)$   $\triangleright$  Smallest  $n_x$  values and indices of the posterior mean in  $X$   
4:  $S_e \leftarrow S_o[I_e, :]$ ,  $S_x \leftarrow X[I_x, :]$   $\triangleright$  Starting points: exploration set and exploitation set  
5:  $[\tilde{\mathbf{x}}^*, \tilde{f}^*] \leftarrow \text{minimize}(f(\mathbf{x}), S)$ ,  $S = S_e \cup S_x$   $\triangleright$  Gradient-based multistart optimization  
**Output:** Thompson sample  $\tilde{\mathbf{x}}^*$   $\triangleright$  Global minimum of the posterior sample

be the  $n_o$  smallest local minima of the prior sample. The set of starting points,  $S$ , is defined as:

$$S = S_e \cup S_x, \quad S_e = \arg \text{mink}_{\mathbf{x} \in S_o}(f_{\tilde{\omega}}(\mathbf{x}), n_e), \quad S_x = \arg \text{mink}_{\mathbf{x} \in X}(f_{\tilde{\omega}}(\mathbf{x}), n_x). \quad (4)$$

Algorithm 1 outlines the procedure for TS-roots. Here,  $n_e$  and  $n_x$  are imposed to control the cost of gradient-based multistart optimization, and  $n_o$  is set to limit the number of evaluations of  $\tilde{f}_{\tilde{\omega}}$  in the determination of  $S_e$ . Considering the cost difference, we can have  $n_o \gg n_e$ . We observe that  $n_e$  and  $n_x$  can be set to small values, and  $n_o$  to a moderate value, without sacrificing the quality of optimization, see Appendix D. The TS-roots algorithm scales linearly in  $d$ , see Appendix C.5.

### 3.2 RELATIONS BETWEEN THE LOCAL MINIMA OF PRIOR AND POSTERIOR SAMPLES

Figure 1 shows several posterior samples  $\tilde{f}_{\tilde{\omega}}$  in one and two dimensions, each marked with its local minima  $\check{\mathbf{x}}_{\tilde{\omega}}$  and global minimum  $\tilde{\mathbf{x}}_{\tilde{\omega}}^*$ . Here the exploration set  $S_e$  is the local minima  $\check{\mathbf{x}}_{\tilde{\omega}}$  of  $f_{\tilde{\omega}}$ , and the exploitation set  $S_x$  is the observed locations  $X$ . We make the following observations:

- (1) The prior sample  $f_{\omega}$  is more complex than the data adjustment  $b$  in the sense that it is less smooth and has more critical points. The comparison of smoothness can be made rigorous in various ways: for example, for GPs with a Matérn covariance function where the smoothness parameter is finite,  $f_{\omega}$  is almost everywhere one time less differentiable than  $b$  (see e.g., Garnett (2023) Sec. 10.2, Kanagawa et al. (2018)).
- (2) Item (1) implies that when the prior sample  $f_{\omega}$  is added to the smoother landscape of  $b$ , each local minimum  $\check{\mathbf{x}}_{\omega}$  of  $f_{\omega}$  will be located near a local minimum  $\check{\mathbf{x}}_{\tilde{\omega}}$  of the posterior sample  $\tilde{f}_{\tilde{\omega}}$ . Away from the observed locations  $X$ , each  $\check{\mathbf{x}}_{\tilde{\omega}}$  is closely associated with an  $\check{\mathbf{x}}_{\omega}$ , with minimal change in location. In the vicinity of  $X$ , an  $\check{\mathbf{x}}_{\tilde{\omega}}$  may have both a data point  $\mathbf{x}^i$  and an  $\check{\mathbf{x}}_{\omega}$  nearby, but because of the smoothness difference of  $f_{\omega}$  and  $b$ , in most cases the one closest to  $\check{\mathbf{x}}_{\tilde{\omega}}$  is an  $\check{\mathbf{x}}_{\omega}$ .
- (3) It is possible that near  $X$ , sharp changes in  $f_{\omega}$  may require sharp adjustments to the data, which may move some  $\check{\mathbf{x}}_{\omega}$  by a significant distance, or create new  $\check{\mathbf{x}}_{\tilde{\omega}}$  that are not near any  $\check{\mathbf{x}}_{\omega}$  or any  $\mathbf{x}^i$ .
- (4) Searching from  $\mathbf{x}^i$  with good observed values can discover good  $\check{\mathbf{x}}_{\tilde{\omega}}$  in the vicinity of  $X$ , which can pick up some local minima not readily discovered by  $\check{X}_{\omega}$ . This is especially true if  $f_{\omega}$  is relatively flat near  $\mathbf{x}^i$ .
- (5) Since the posterior sample  $\tilde{f}_{\tilde{\omega}}$  adapts to the dataset, searches from  $\mathbf{x}^i$  will tend to converge to a good  $\check{\mathbf{x}}_{\tilde{\omega}}$  among all the local optima near  $\mathbf{x}^i$ . Even if the searches from  $X$  cannot discover all the local minima in its vicinity, they tend to discover a good subset of them. Therefore, (4) can help address the issue in (3), if not fully eliminating it. By combining subsets of  $\check{X}_{\omega}$  and  $X$ , we can expect that the set of local minima of  $\tilde{f}_{\tilde{\omega}}$  discovered with these starting points includes the global minimum  $\tilde{\mathbf{x}}_{\tilde{\omega}}^*$  with a high probability with respect to  $\tilde{\omega}$ .

### 3.3 A REPRESENTATION OF PRIOR SAMPLE LOCAL MINIMA

For each component function  $f_i(x_i; \omega_i)$  of the prior sample  $f_{\omega}(\mathbf{x})$ , define  $h_i(\underline{x}_i) = f'_i(\underline{x}_i)$ ,  $h_i(\bar{x}_i) = -f'_i(\bar{x}_i)$ , and  $h_i(x_i) = f''_i(x_i)$  for  $x_i \in (\underline{x}_i, \bar{x}_i)$ . We call a coordinate  $\xi_i \in [\underline{x}_i, \bar{x}_i]$  of *mono type* if  $f_i(\xi_i)h_i(\xi_i) > 0$  and call it of *mixed type* if  $f_i(\xi_i)h_i(\xi_i) < 0$ . Let  $\tilde{\Xi}_i = \{\xi_{i,j}\}_{j=1}^{r_i}$  be the set

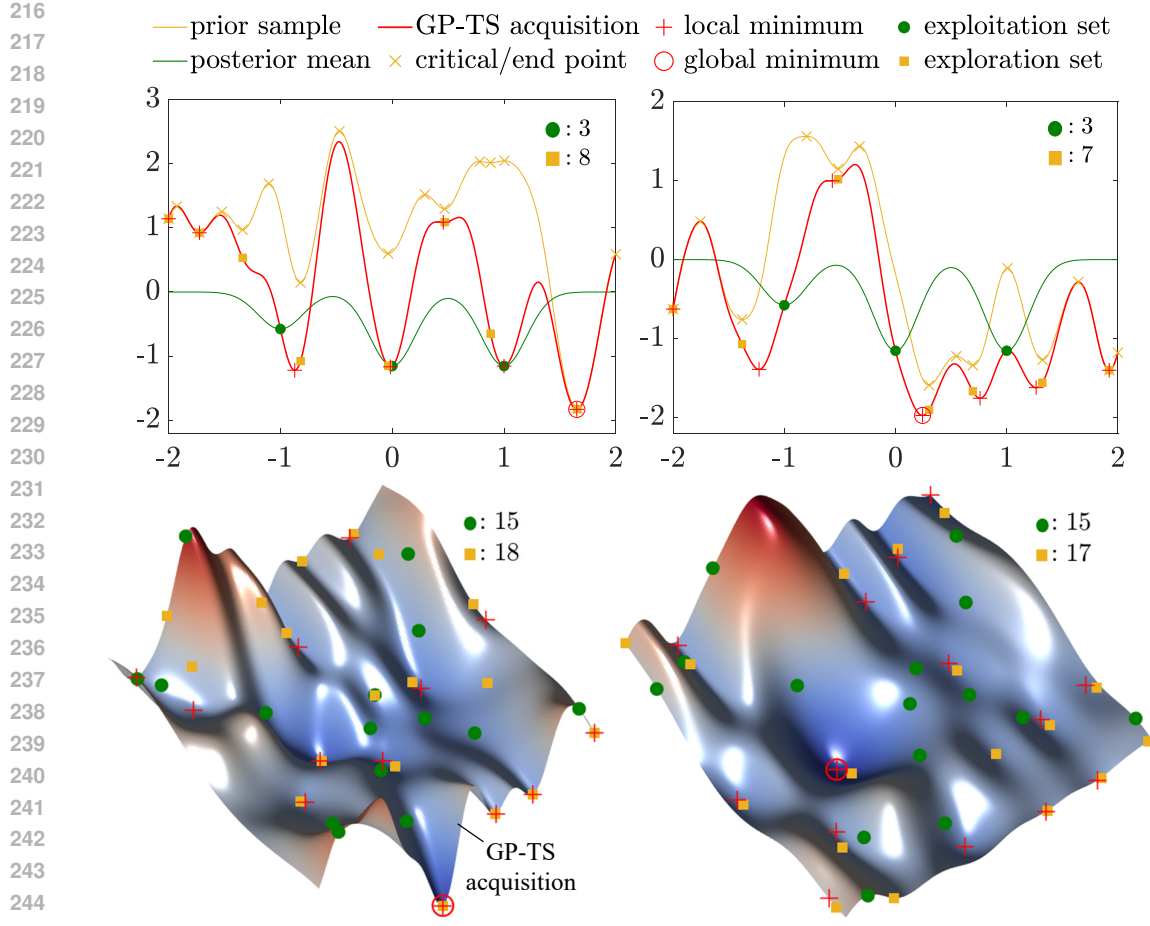


Figure 1: Illustrations of exploration and exploitation sets for the global optimization of GP-TS acquisition functions in one dimension (top row) and two dimensions (bottom row). *Left column:* When the global minimum  $\tilde{\mathbf{x}}_{\omega}^*$  of the GP-TS acquisition function lies outside the interpolation region, it is typically identified by starting the gradient-based multistart optimizer at a local minimum of the prior sample. *Right column:* When  $\tilde{\mathbf{x}}_{\omega}^*$  is within the interpolation region, it can be found by starting the optimizer at either an observed location or a local minimum of the prior sample.

of interior critical points of  $f_i$  such that  $\xi_{i,j} \in (\underline{x}_i, \bar{x}_i)$  and  $f'_i(\xi_{i,j}) = 0$ ,  $j \in \{1, \dots, r_i\}$ . Denote  $\xi_{i,0} = \underline{x}_i$  and  $\xi_{i,r_i+1} = \bar{x}_i$ . Partition the set of candidate coordinates  $\Xi_i = \{\xi_{i,j}\}_{j=0}^{r_i+1}$  into mono type  $\Xi_i^{(0)}$  and mixed type  $\Xi_i^{(1)}$ . Proposition 1 gives a representation of the sets of strong local extrema of the prior sample. [Its proof and the set sizes therein are given in Appendix A.](#)

**Proposition 1** *The set of strong local minima of the prior sample  $f_{\omega}(\mathbf{x})$  can be written as:*

$$\check{X}_{\omega} = \check{X}_{\omega}^{-} \sqcup \check{X}_{\omega}^{+}, \quad \check{X}_{\omega}^{-} = \{\boldsymbol{\xi} \in \Xi^{(1)} : f_{\omega}(\boldsymbol{\xi}) < 0\}, \quad \check{X}_{\omega}^{+} = \{\boldsymbol{\xi} \in \Xi^{(0)} : f_{\omega}(\boldsymbol{\xi}) > 0\}, \quad (5)$$

where tensor grids  $\Xi^{(j)} = \prod_{i=1}^d \Xi_i^{(j)}$ ,  $j \in \{0, 1\}$ . The set  $\hat{X}_{\omega}$  of strong local maxima of  $f_{\omega}(\mathbf{x})$  has an analogous representation, and satisfies  $\hat{X}_{\omega} \sqcup \check{X}_{\omega} = \Xi^{(0)} \sqcup \Xi^{(1)}$ , where  $\sqcup$  is the disjoint union.

**Critical Points of Univariate Functions via Global Rootfinding.** To compute the set  $\Xi_i$  of all critical points of  $f_i$  is to compute all the roots of the derivative  $f'_i$  on the interval  $[\underline{x}_i, \bar{x}_i]$ . Since  $f'_i$  is continuous, this can be solved robustly and efficiently by approximating the function with a Chebyshev or Legendre polynomial and solving a structured eigenvalue problem (see e.g., Trefethen (2019)). The `roots` algorithm for global rootfinding based on polynomial approximation is given as Algorithm 3 in Appendix C.

### 3.4 ORDERING OF PRIOR SAMPLE LOCAL MINIMA

While the size of  $\check{X}_\omega$  grows exponentially in domain dimension  $d$ , its representation in eq. (5) enables an efficient algorithm to compute the best subset  $S_o$  (eq. (3)) without enumerating its elements.

With eq. (5), we see that  $\check{X}_\omega^-$  consists of all the local minima of  $f_\omega$  with negative function values. Consider the case where  $\check{X}_\omega^-$  has at least  $n_o$  elements so that in the definition of  $S_o$  we can replace  $\check{X}_\omega$  with  $\check{X}_\omega^-$ , which in turn can be replaced with  $\Xi^{(1)}$ . As we will show later, the problem of finding the largest elements of  $|f_\omega(\mathbf{x})|$  within  $\Xi^{(1)}$  is easier than finding the smallest negative elements of  $f_\omega(\mathbf{x})$ . Once the former is solved, we can solve the latter simply by removing the positive elements. Therefore, we convert the problem of eq. (3) into three steps:

1.  $S^{(1)} = \arg \max_{\mathbf{x} \in \Xi^{(1)}} (|f_\omega(\mathbf{x})|, \alpha n_o)$ , with buffer coefficient  $\alpha \geq 1$ ;
2.  $\check{S}^- = \{\mathbf{x} \in S^{(1)} : f_\omega(\mathbf{x}) < 0\}$ , so that  $\check{S}^- \subseteq \check{X}_\omega^-$ ;
3.  $S_o = \arg \min_{\mathbf{x} \in \check{S}^-} (f_\omega(\mathbf{x}), n_o)$ , assuming that  $|\check{S}^-| \geq n_o$ .

The last two steps are by enumeration and straightforward. The first step can be solved efficiently using min-heaps, with a time complexity that scales linearly in  $\sum_{i=1}^d |\Xi_i^{(1)}|$  rather than  $\prod_{i=1}^d |\Xi_i^{(1)}|$ , see Appendix B. The coefficient  $\alpha$  is chosen so that  $|\check{S}^-| \geq n_o$ . The case when  $|\check{X}_\omega^-| < n_o < |\check{X}_\omega|$  can be handled similarly. If  $n_o \geq |\check{X}_\omega|$ , no subsetting is needed. The overall procedure to compute  $S_o$  is given in Algorithm 4 in Appendix C.

## 4 SAMPLE-AVERAGE POSTERIOR FUNCTION

We finally propose a sample-average posterior function that explicitly controls the exploration–exploitation balance and, notably, can be generated at the cost of generating one posterior sample. Let  $\tilde{\mu}(\mathbf{x}) = \kappa_{\cdot, n}(\mathbf{x})(\mathbf{K}_{n, n} + \Sigma)^{-1} \mathbf{y}$  be the posterior mean function. For noiseless observations with  $\tilde{\omega} = \omega$ , we can rewrite the GP posterior function in eq. (2) as  $\tilde{f}_\omega(\mathbf{x}) = f_\omega(\mathbf{x}) + \tilde{\mu}(\mathbf{x}) + \xi_\omega(\mathbf{x})$ , where  $\xi_\omega(\mathbf{x}) = -\kappa_{\cdot, n}(\mathbf{x})\mathbf{K}_{n, n}^{-1}\mathbf{f}_n$ . Define  $\alpha_{\text{aTS}}(\mathbf{x}) = \frac{1}{N_c} \sum_{j=1}^{N_c} \tilde{f}_\omega^j(\mathbf{x})$  as the sample-average posterior function, where  $\tilde{f}_\omega^j(\mathbf{x})$  are samples generated from the GP posterior and  $N_c \in \mathbb{N}_{>0}$ . Since  $\tilde{\mu}(\mathbf{x})$  is deterministic, and the scaled prior sample  $\frac{1}{\sqrt{N_c}} f_\omega^j(\mathbf{x})$  can be written as  $\frac{1}{\sqrt{N_c}} f_\omega^j(\mathbf{x}) \stackrel{\text{iid}}{\sim} \mathcal{GP}(0, \frac{1}{N_c} \kappa(\mathbf{x}, \mathbf{x}'))$ , we have  $\alpha_{\text{aTS}}(\mathbf{x}) = \mu(\mathbf{x}) + \frac{1}{\sqrt{N_c}} (f_\omega(\mathbf{x}) + \xi_\omega(\mathbf{x}))$ , where the first and second terms favor exploitation and exploration, respectively. Thus, we can consider  $N_c$  as an exploration–exploitation control parameter that, at large values, prioritizes exploitation by concentrating the conditional distribution of the global minimum location, i.e.,  $p(\mathbf{x}^* | \mathcal{D})$ , at the minimum location of  $\tilde{\mu}(\mathbf{x})$ , see Figure 12 in Appendix I. With  $\alpha_{\text{aTS}}(\mathbf{x})$ , we can reproduce  $f_\omega(\mathbf{x})$  and the GP mean function  $\tilde{\mu}(\mathbf{x})$  by setting  $N_c = 1$  and  $N_c = \infty$ , respectively.

## 5 RELATED WORKS

**Sampling from Gaussian Process Posteriors.** A prevalent method to sample GP posteriors with stationary covariance functions is via weight-space approximations based on Bayesian linear models of random Fourier features (Rahimi & Recht, 2007). This method, unfortunately, is subject to the variance starvation problem (Mutny & Krause, 2018; Wilson et al., 2020) which can be mitigated using more accurate feature representations (see e.g., Hensman et al. (2018); Solin & Särkkä (2020)). An alternative is pathwise conditioning (Wilson et al., 2020) that draws GP posterior samples by updating the corresponding prior samples. The decoupled representation of the pathwise conditioning can be further reformulated as two stochastic optimization problems for the posterior mean and an uncertainty reduction term, which are then efficiently solved using stochastic gradient descent (Lin et al., 2023).

**Optimization of Acquisition Functions.** While their global optima guarantee the Bayes’ decision rule, BO acquisition functions are highly non-convex and difficult to optimize (Wilson et al., 2018). Nevertheless, less attention has been given to the development of robust algorithms for optimizing

these acquisition functions. For this inner-loop optimization, gradient-based optimizers are often selected because of their fast convergence and robust performance (Daulton et al., 2020). The implementation of such optimizers is facilitated by Monte Carlo (MC) acquisition functions whose derivatives are easy to evaluate (Wilson et al., 2018). Gradient-based optimizers also allow multistart settings that use a set of starting points which can be, for example, midpoints of data points (Jones, 2001), uniformly distributed samples over the input variable space (Frazier, 2018; Ament et al., 2023), or random points from a Latin hypercube design (Wang et al., 2020). However, multistart-based methods with random search may have difficulty determining the non-flat regions of acquisition functions, especially in high dimensions (Rana et al., 2017). [The log reformulation approach is a good solution to the numerical pathology of flat acquisition surface over large regions of the input variable space \(Ament et al., 2023\). While this approach works for acquisition functions prone to the flat surface issue such as the family of EI-based acquisition functions, its performance has yet to be evaluated for acquisition functions with many local minima like those based on posterior samples.](#)

**Posterior Sample-Based Acquisition Functions.** As discussed in Section 1, the family of posterior sample-based acquisition functions is determined from samples of the posterior. GP-TS (Chowdhury & Gopalan, 2017) is a notable member that extends the classical TS for finite-armed bandits to continuous settings of BO (see algorithms in Appendix E). GP-TS prefers exploration by the mechanism that iteratively samples a function from the GP posterior of the objective function, optimizes this function, and selects the resulting solution as the next candidate for objective evaluation. To further improve the exploitation of GP-TS, the sample mean of MC acquisition functions can be defined from multiple samples of the posterior (Wilson et al., 2018; Balandat et al., 2020). Different types of MC acquisition functions can also be used to inject beliefs about functions into the prior (Hvarfner et al., 2024).

## 6 RESULTS

We assess the performance of TS-roots in optimizing benchmark functions. We then compare the quality of solutions to the inner-loop optimization of GP-TS acquisition functions obtained from our proposed method, a gradient-based multistart optimizer with uniformly random starting points, and a genetic algorithm. We also show how TS-roots can improve the performance of MES. Finally, we propose a new sample-average posterior function and show how it affects the performance of GP-TS. The experimental details for the presented results are in Appendix H.

**Optimizing Benchmark Functions.** We test the empirical performance of TS-roots on challenging minimization problems of five benchmark functions: the 2D Schwefel, 4D Rosenbrock, 10D Levy, 16D Ackley, and 16D Powell functions (Surjanovic & Bingham, 2013). The analytical expressions for these functions and their global minimum are given in Appendix F.

In each optimization iteration, we record the best observed value of the error  $\log(y_{\min} - f^*)$  and the distance  $\log(\|\mathbf{x}_{\min} - \mathbf{x}^*\|)$ , where  $y_{\min}$ ,  $\mathbf{x}_{\min}$ ,  $f^*$ , and  $\mathbf{x}^*$  are the best observation of the objective function in each iteration, the corresponding location of the observation, the true minimum value of the objective function, and the true minimum location, respectively. We compare the optimization results obtained from TS-roots and other BO methods, including GP-TS using decoupling sampling with random Fourier features (TS-DSRF), GP-TS with random Fourier features (TS-RF), expected improvement (EI) (Jones et al., 1998), and lower confidence bound (LCB)—the version of GP-UCB (Srinivas et al., 2010) for minimization problems.

Figure 2 shows the medians and interquartile ranges of solution values obtained from 20 runs of each of the considered BO methods. The corresponding histories of solution locations are in Figure 9 of Appendix I. With a fair comparison of outer-loop results (detailed in Appendix H), TS-roots surprisingly shows the best performance on the 2D Schwefel, 16D Ackley, and 16D Powell functions, and gives competitive results in the 4D Rosenbrock and 10D Levy problems. Notably, TS-roots recommends better solutions than its counterparts, TS-DSRF and TS-RF, in high-dimensional problems and offers competitive performance in low-dimensional problems. Across all the examples, EI and LCB tend to perform better in the initial stages, while TS-roots shows fast improvement in later stages. This is because GP-TS favors exploration, which delays rewards. The exploration phase, in general, takes longer for higher-dimensional problems.

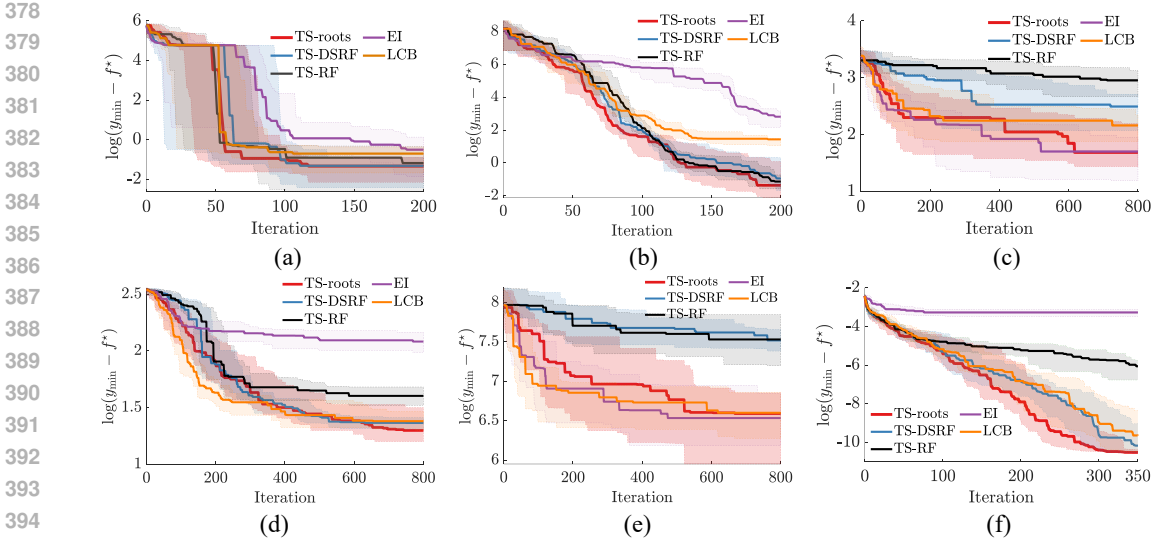


Figure 2: Outer-loop optimization results for the (a) 2D Schwefel function, (b) 4D Rosenbrock function, (c) 10D Levy function, (d) 16D Ackley function, (e) 16D Powell function, and (f) ten-bar truss problem. The plots are histories of medians and interquartile ranges of solution values from 20 runs of TS-roots, TS-DSRF (i.e., TS using decoupled sampling with random Fourier features), TS-RF (i.e., TS using random Fourier features), EI, and LCB.

**Optimizing Real-world Problem.** We implement TS-roots to optimize an engineered ten-bar truss structure (see Appendix G). Ten design variables of the truss are the cross-sectional areas of the truss members. The objective is to minimize a weighted sum of the scaled total cross-sectional area and the scaled vertical displacement at a node of interest.

Figure 2(f) shows the outer-loop optimization results for the truss obtained from 20 runs of each BO method, where  $f^*$  is a lower bound of the best objective function value we observed from all runs. TS-roots provides the best optimization result with rapid convergence.

**Optimizing GP-TS Acquisition Functions via Rootfinding.** We assess the quality of solutions and computational cost for the inner-loop optimization of GP-TS acquisition functions by the proposed global optimization algorithm, referred to as rootfinding hereafter. We do so by computing the optimized values  $\alpha_k^*$  of the GP-TS acquisition functions, the corresponding solution points  $\mathbf{x}_k^*$ , and the CPU times  $t_k$  required for optimizing the acquisition functions during the optimization process. For low-dimensional problems of the 2D Schwefel and 4D Rosenbrock functions, we also compute the exact global solution points  $\mathbf{x}_k^t$  of the GP-TS acquisition functions by starting the gradient-based optimizer at a large number of initial points (set as  $10^4$ ), which is much larger than the maximum number of starting points set for TS-roots. For comparison, we extend the same GP-TS acquisition functions to inner-loop optimization using a gradient-based multistart optimizer with random starting points (i.e., random multistart) and a genetic algorithm. In each outer-loop optimization iteration, the number of starting points for the random multistart and the population size of the genetic algorithm are equal to the number of starting points recommended for rootfinding. The same termination conditions are used for the three algorithms.

Figure 3 shows the comparative performance of the inner-loop optimization for low-dimensional cases: the 2D Schwefel and 4D Rosenbrock functions. We see that the optimized acquisition function values and the optimization runtimes by rootfinding and the random multistart algorithm are almost identical, both of which are much better than those by the genetic algorithm. Rootfinding gives the best quality of the new solution points in both cases, while the genetic algorithm gives the worst. Notably in higher-dimensional settings of the 10D Levy, 16D Ackley, and 16D Powell functions shown in Figure 4, rootfinding performs much better than the random multistart and genetic algorithm in terms of optimized acquisition values and optimization runtimes, which verifies the importance of the judicious selection of starting points for global optimization of the GP-TS acquisition functions and the efficiency of rootfinding in high dimensions. The performance of



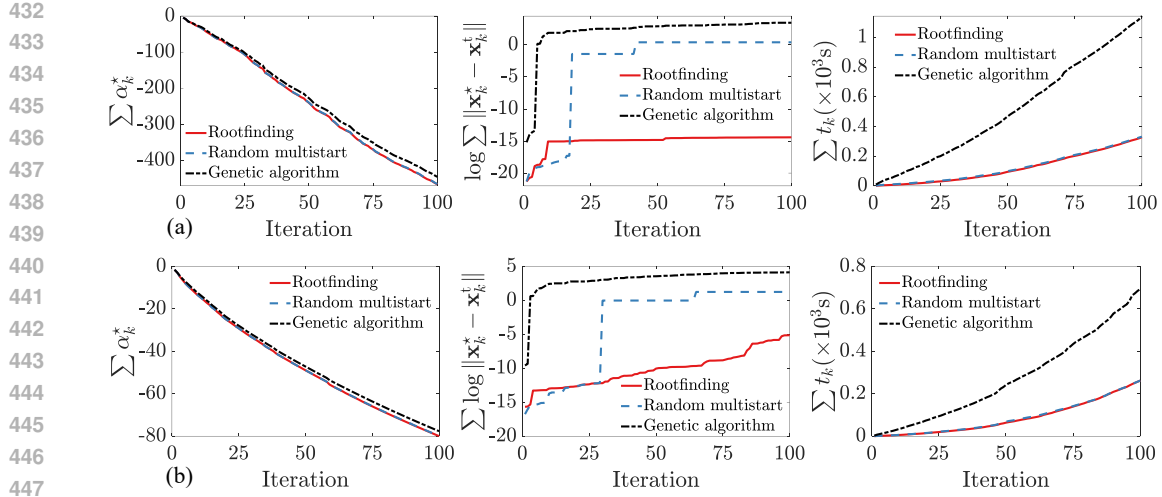


Figure 3: Inner-loop optimization results by rootfinding, a gradient-based multistart optimizer with random starting points (random multistart), and a genetic algorithm for (a) the 2D Schwefel and (b) 4D Rosenbrock functions. The plots are cumulative values of optimized GP-TS acquisition functions  $\alpha_k^*$ , cumulative distances between new solution points  $\mathbf{x}_k^*$  and the true global minima  $\mathbf{x}_k^{\dagger}$  of the acquisition functions, and cumulative CPU times  $t_k$  for optimizing the acquisition functions.

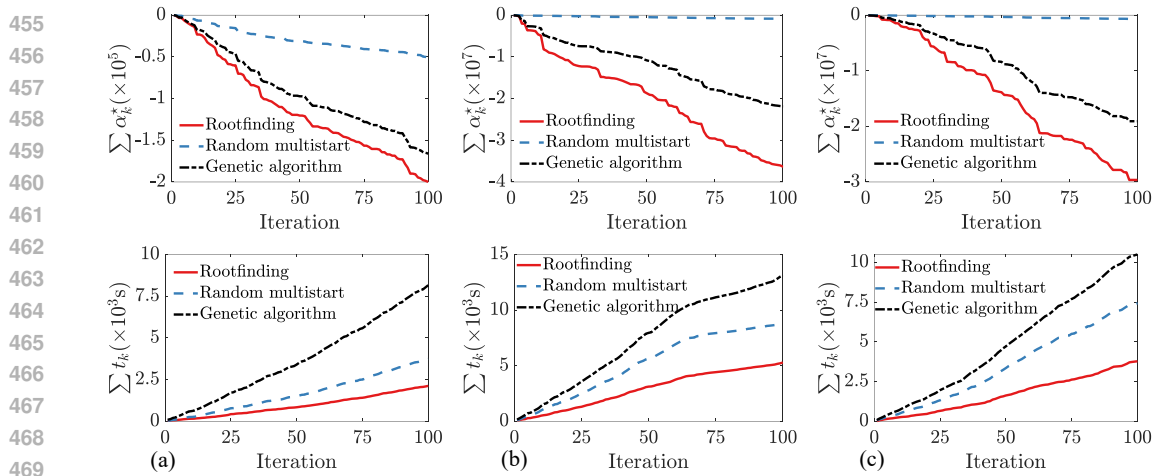


Figure 4: Inner-loop optimization results by rootfinding, a gradient-based multistart optimizer with random starting points (random multistart), and a genetic algorithm for (a) the 10D Levy, (b) 16D Ackley, and (c) 16D Powell functions. The plots are cumulative values of optimized GP-TS acquisition functions  $\alpha_k^*$  and cumulative CPU times  $t_k$  for optimizing the acquisition functions.

random multistart becomes worse in higher dimensions. [Appendix I](#) provides additional results for gradient-based multistart optimization using two other initialization schemes: uniform grid and Latin hypercube sampling. Rootfinding outperforms both, especially in higher dimensions.

**TS-roots to Information-Theoretic Acquisition Functions.** We show how TS-roots can enhance the performance of MES (Wang & Jegelka, 2017), which uses information about the maximum function value  $f^*$  for conducting BO. One approach to computing MES generates a set of GP posterior samples using TS-RF and subsequently optimizes the generated functions for samples of  $f^*$  using a gradient-based multistart optimizer with a large number of generated random starting points (Wang & Jegelka, 2017). We hypothesize that high-quality  $f^*$  samples can improve the performance of MES. Thus, we assign both TS-roots and TS-RF as the inner workings of MES and then compare the result-

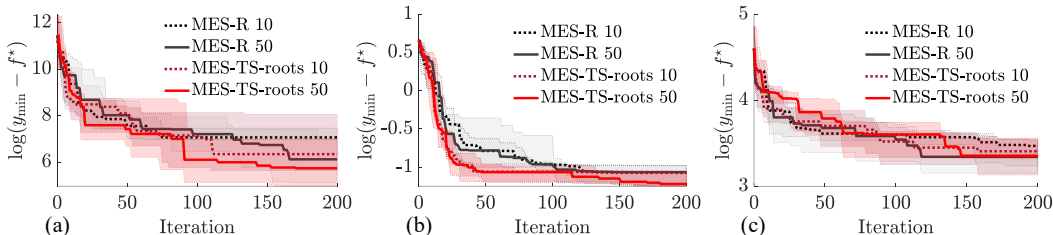


Figure 5: Performance of MES-R 10 and MES-R 50 for (a) the 4D Rosenbrock function, (b) the 6D Hartmann function, and (c) the 10D Levy function when TS-RF and TS-roots are used for generating random samples from  $f^*|\mathcal{D}$ . The plots are histories of medians and interquartile ranges of solutions from ten runs of each method.

ing optimal solutions. Note that the inner-loop optimization of MES, which strongly influences the optimization results, is not addressed by TS-roots.

Specifically, we minimize the 4D Rosenbrock, 6D Hartmann, and 10D Levy functions using four versions of MES, namely MES-R 10, MES-R 50, MES-TS-roots 10, and MES-TS-roots 50. Here, MES-R (Wang & Jegelka, 2017) and MES-TS-roots correspond to TS-RF and TS-roots, respectively, while 10 and 50 represent the number of random samples  $f^*$  generated for computing the MES acquisition function in each iteration.

Figure 5 shows the optimization histories for ten independent runs of each MES method. On the 4D Rosenbrock and 6D Hartmann functions, MES with TS-roots demonstrates superior optimization performance and faster convergence compared to MES with TS-RF, especially when 50 samples of  $f^*$  are generated. For the 10D Levy function, TS-roots outperforms TS-RF when using 10 samples of  $f^*$ , while their performance is comparable when 50 samples are used.

**Performance of Sample-Average Posterior Functions.** We investigate how  $\alpha_{\text{aTS}}(\mathbf{x})$  influences the outer-loop optimization results. For this, we set  $N_c \in \{1, 10, 50, 100\}$  for TS-roots to optimize the 2D Schwefel, 4D Rosenbrock, and 6D Ackley functions. We observe that increasing  $N_c$  from 1 to 10 improves TS-roots performance on the 2D Schwefel, 4D Rosenbrock, and 6D Ackley functions (see Figure 13 in Appendix I). However, further increases in  $N_c$  from 10 to 50 and 100 result in slight declines in solution quality as TS-roots transitions to exploitation. These observations indicate that there is an optimal value of  $N_c$  for each problem at which TS-roots achieves its best performance by balancing exploitation and exploration. However, identifying the optimal value to maximize the performance of  $\alpha_{\text{aTS}}(\mathbf{x})$  for a particular optimization problem remains an open issue.

## 7 CONCLUSION AND FUTURE WORK

We presented TS-roots, a global optimization strategy for posterior sample paths. It features an adaptive selection of starting points for gradient-based multistart optimizers, combining exploration and exploitation. This strategy breaks the curse of dimensionality by exploiting the separability of Gaussian process priors. Compared with random multistart and a genetic algorithm, TS-roots consistently yields higher-quality solutions in optimizing posterior sample paths, across a range of input dimensions. It also improves the outer-loop optimization performance of GP-TS and information-theoretic acquisition functions such as MES for Bayesian optimization. For future work, we aim to extend TS-roots to other spectral representations per Bochner’s theorem (Mutny & Krause, 2018; Hensman et al., 2018; Solin & Särkkä, 2020). We also plan to study the ways and the probability of TS-roots failing to find the global optimum, as well as the impact of subset sizes.

## REFERENCES

- 540  
541  
542 Sebastian Ament, Samuel Daulton, David Eriksson, Maximilian Balandat, and Eytan Bak-  
543 shy. Unexpected improvements to expected improvement for Bayesian optimization.  
544 In *Advances in Neural Information Processing Systems*, volume 36, pp. 20577–20612,  
545 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/  
546 file/419f72cbd568ad62183f8132a3605a2a-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/419f72cbd568ad62183f8132a3605a2a-Paper-Conference.pdf).
- 547 Maximilian Balandat, Brian Karrer, Daniel Jiang, Samuel Daulton, Ben Letham, Andrew G Wil-  
548 son, and Eytan Bakshy. BoTorch: A framework for efficient Monte-Carlo Bayesian opti-  
549 mization. In *Advances in Neural Information Processing Systems*, volume 33, pp. 21524–  
550 21538, 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
551 2020/file/f5b1b89d98b7286673128a5fb112cb9a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f5b1b89d98b7286673128a5fb112cb9a-Paper.pdf).
- 552 Zachary Battles and Lloyd N. Trefethen. An extension of MATLAB to continuous functions  
553 and operators. *SIAM Journal on Scientific Computing*, 25(5):1743–1770, 2004. doi: 10.1137/  
554 S1064827503430126.
- 555 James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal*  
556 *of Machine Learning Research*, 13(10):281–305, 2012. URL [http://jmlr.org/papers/  
557 v13/bergstral2a.html](http://jmlr.org/papers/v13/bergstral2a.html).
- 558 Adam D. Bull. Convergence rates of efficient global optimization algorithms. *Journal of Ma-*  
559 *chine Learning Research*, 12(88):2879–2904, 2011. URL [http://jmlr.org/papers/  
560 v12/bull11a.html](http://jmlr.org/papers/v12/bull11a.html).
- 561 Olivier Chapelle and Lihong Li. An empirical evaluation of Thompson sampling. In  
562 *Advances in Neural Information Processing Systems*, volume 24, pp. 2249–2257,  
563 2011. URL [https://papers.nips.cc/paper\\_files/paper/2011/hash/  
564 e53a0a2978c28872a4505bdb51db06dc-Abstract.html](https://papers.nips.cc/paper_files/paper/2011/hash/e53a0a2978c28872a4505bdb51db06dc-Abstract.html).
- 565 Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *Proceedings*  
566 *of the 34th International Conference on Machine Learning*, volume 70, pp. 844–853, 06–11 Aug  
567 2017. URL <https://proceedings.mlr.press/v70/chowdhury17a.html>.
- 568 Samuel Daulton, Maximilian Balandat, and Eytan Bakshy. Differentiable expected hyper-  
569 volume improvement for parallel multi-objective bayesian optimization. In *Advances*  
570 *in Neural Information Processing Systems*, volume 33, pp. 9851–9864, 2020. URL  
571 [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/  
572 6fec24eac8f18ed793f5eaad3dd7977c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/6fec24eac8f18ed793f5eaad3dd7977c-Paper.pdf).
- 573 David Duvenaud, James Lloyd, Roger Grosse, Joshua Tenenbaum, and Ghahramani Zoubin. Struc-  
574 ture discovery in nonparametric regression through compositional kernel search. In *Proceedings*  
575 *of the 30th International Conference on Machine Learning*, volume 28, pp. 1166–1174, 2013.  
576 URL <http://proceedings.mlr.press/v28/duvenaud13.html>.
- 577 Thomas A. Feo and Mauricio G. C. Resende. Greedy randomized adaptive search procedures.  
578 *Journal of Global Optimization*, 6(2):109–133, 1995. doi: 10.1007/BF01096763. URL <https://doi.org/10.1007/BF01096763>.
- 579 Peter I. Frazier. Bayesian optimization. In *Recent Advances in Optimization and Modeling of*  
580 *Contemporary Problems*, INFORMS TutORials in Operations Research, chapter 11, pp. 255–  
581 278. October 2018. doi: 10.1287/educ.2018.0188.
- 582 Roman Garnett. *Bayesian Optimization*. Cambridge University Press, Cambridge, UK, 2023. doi:  
583 10.1017/9781108348973.
- 584 I.S. Gradshteyn and I.M. Ryzhik. *Table of Integrals, Series, and Products*. Academic Press, Boston,  
585 8th edition, 2014. ISBN 978-0-12-384933-5. doi: 10.1016/C2010-0-64839-5.
- 586 Nikolaus Hansen, Sibylle D. Müller, and Petros Koumoutsakos. Reducing the time complexity of  
587 the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). *Evolutionary*  
588 *Computation*, 11(1):1–18, 2003. doi: 10.1162/106365603321828970.

- 594 Philipp Hennig and Christian J. Schuler. Entropy search for information-efficient global opti-  
595 mization. *Journal of Machine Learning Research*, 13(6):1809–1837, 2012. URL <https://www.jmlr.org/papers/v13/hennig12a.html>.  
596  
597
- 598 James Hensman, Nicolas Durrande, and Arno Solin. Variational Fourier features for Gaussian pro-  
599 cesses. *Journal of Machine Learning Research*, 18(151):1–52, 2018. URL <http://jmlr.org/papers/v18/16-579.html>.  
600
- 601 José Miguel Hernández-Lobato, James Requeima, Edward O. Pyzer-Knapp, and Alán Aspuru-  
602 Guzik. Parallel and distributed Thompson sampling for large-scale accelerated exploration of  
603 chemical space. In *Proceedings of the 34th International Conference on Machine Learning*, vol-  
604 ume 70, pp. 1470–1479, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/hernandez-lobato17a.html>.  
605  
606
- 607 José Miguel Hernández-Lobato, Matthew W. Hoffman, and Zoubin Ghahramani. Predic-  
608 tive entropy search for efficient global optimization of black-box functions. In *Ad-  
609 vances in Neural Information Processing Systems*, volume 27, pp. 918–926, 2014. URL  
610 [https://proceedings.neurips.cc/paper\\_files/paper/2014/hash/  
611 069d3bb002acd8d7dd095917f9efe4cb-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2014/hash/069d3bb002acd8d7dd095917f9efe4cb-Abstract.html).
- 612 Yehuda Hoffman and Erez Ribak. Constrained Realizations of Gaussian Fields: A Simple Algo-  
613 rithm. *Astrophysical Journal Letters*, 380:L5–L8, October 1991. doi: 10.1086/186160. URL  
614 <https://ui.adsabs.harvard.edu/abs/1991ApJ...380L...5H>.  
615
- 616 Carl Hvarfner, Frank Hutter, and Luigi Nardi. Joint entropy search for  
617 maximally-informed Bayesian optimization. In *Advances in Neural Infor-  
618 mation Processing Systems*, volume 35, pp. 11494–11506, 2022. URL  
619 [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/  
620 4b03821747e89ce803b2dac590f6a39b-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/4b03821747e89ce803b2dac590f6a39b-Abstract-Conference.html).
- 621 Carl Hvarfner, Frank Hutter, and Luigi Nardi. A general framework for user-guided bayesian opti-  
622 mization. In *The Twelfth International Conference on Learning Representations*, pp. 9851–9864,  
623 2024. URL <https://iclr.cc/virtual/2024/poster/18774>.  
624
- 625 D. R. Jones, C. D. Perttunen, and B. E. Stuckman. Lipschitzian optimization without the Lipschitz  
626 constant. *Journal of Optimization Theory and Applications*, 79(1):157–181, 1993. doi: 10.1007/  
627 BF00941892. URL <https://doi.org/10.1007/BF00941892>.  
628
- 629 Donald R. Jones. A taxonomy of global optimization methods based on response surfaces. *Journal  
630 of Global Optimization*, 21(4):345–383, 2001. doi: 10.1023/A:1012771025575.
- 631 Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of  
632 expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998. doi:  
633 10.1023/A:1008306431147. URL <https://doi.org/10.1023/A:1008306431147>.  
634
- 635 Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian  
636 processes and kernel methods: A review on connections and equivalences, 2018. URL <https://arxiv.org/abs/1807.02582>.  
637
- 638 Kirthevasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabas Poczos. Parallelised  
639 Bayesian optimisation via Thompson sampling. In *Proceedings of the Twenty-First International  
640 Conference on Artificial Intelligence and Statistics*, volume 84, pp. 133–142, 09–11 Apr 2018.  
641 URL <https://proceedings.mlr.press/v84/kandasamy18a.html>.  
642
- 643 Jungtaek Kim and Seungjin Choi. On local optimizers of acquisition functions in Bayesian opti-  
644 mization. In *Machine Learning and Knowledge Discovery in Databases*, pp. 675–690, 2021. doi:  
645 10.1007/978-3-030-67661-2\_40.
- 646  
647 Harold J. Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in  
the presence of noise. *Journal Basic Engineering*, 86(1):97–106, 1964. doi: 10.1115/1.3653121.

- 648 Jihao Andreas Lin, Javier Antorán, Shreyas Padhy, David Janz, José Miguel Hernández-Lobato,  
649 and Alexander Terenin. Sampling from Gaussian process posteriors using stochastic gradient  
650 descent. In *Advances in Neural Information Processing Systems*, volume 36, pp. 36886–  
651 36912, 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
652 2023/file/7482e8ce4139df1a2d8195a0746fa713-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/7482e8ce4139df1a2d8195a0746fa713-Paper-Conference.pdf).
- 653 David J. C. MacKay. *Information theory, inference, and learning algorithms*. Cambridge University  
654 Press, Cambridge, UK, 2003. URL <https://www.cambridge.org/9780521642989>.
- 655  
656 Melanie Mitchell. *An introduction to genetic algorithms*. The MIT Press, Massachusetts, USA,  
657 1998.
- 658  
659 Mojmir Mutny and Andreas Krause. Efficient high dimensional Bayesian optimization with addi-  
660 tivity and quadrature Fourier features. In *Advances in Neural Information Processing Systems*,  
661 volume 31, pp. 9005–9016, 2018. URL [https://proceedings.neurips.cc/paper/  
662 2018/hash/4e5046fc8d6a97d18a5f54beaed54dea-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/4e5046fc8d6a97d18a5f54beaed54dea-Abstract.html).
- 663  
664 Art B. Owen. A Central Limit Theorem for Latin Hypercube Sampling. *Journal of the Royal Statisti-  
665 cal Society: Series B (Methodological)*, 54(2):541–551, 12 1992. doi: 10.1111/j.2517-6161.1992.  
666 tb01895.x. URL <https://doi.org/10.1111/j.2517-6161.1992.tb01895.x>.
- 667  
668 Victor Picheny, Tobias Wagner, and David Ginsbourger. A benchmark of kriging-based infill criteria  
669 for noisy optimization. *Structural and multidisciplinary optimization*, 48:607–626, 2013. doi: 10.  
670 1007/s00158-013-0919-4. URL <https://doi.org/10.1007/s00158-013-0919-4>.
- 671  
672 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines.  
673 In *Advances in Neural Information Processing Systems*, volume 20, pp. 1177–1184,  
674 2007. URL [https://proceedings.neurips.cc/paper\\_files/paper/2007/  
675 file/013a006f03dbc5392effeb8f18fda755-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf).
- 676  
677 Santu Rana, Cheng Li, Sunil Gupta, Vu Nguyen, and Svetha Venkatesh. High dimensional  
678 Bayesian optimization with elastic Gaussian process. In *Proceedings of the 34th Interna-  
679 tional Conference on Machine Learning*, volume 70, pp. 2883–2891, 06–11 Aug 2017. URL  
680 <https://proceedings.mlr.press/v70/rana17a.html>.
- 681  
682 Carl Edward Rasmussen and Christopher K I Williams. *Gaussian processes for machine learning*.  
683 The MIT Press, Massachusetts, USA, 2006. ISBN 9780521872508. doi: 10.7551/mitpress/3206.  
684 001.0001. URL <https://doi.org/10.7551/mitpress/3206.001.0001>.
- 685  
686 Mark Richardson. `chebpy`, a Python implementation of `chebfun`, 2016. URL [https://github.  
687 com/chebpy/chebpy](https://github.com/chebpy/chebpy).
- 688  
689 Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of  
690 Operations Research*, 39(4):1221–1243, April 2014. doi: 10.1287/moor.2014.0650.
- 691  
692 Daniel J. Russo and Benjamin Van Roy. An information-theoretic analysis of Thompson sampling.  
693 *The Journal of Machine Learning Research*, 17(1):2442–2471, 2016. URL [https://www.  
694 jmlr.org/papers/v17/14-087.html](https://www.jmlr.org/papers/v17/14-087.html).
- 695  
696 Amar Shah and Zoubin Ghahramani. Parallel predictive entropy search for batch global optimiza-  
697 tion of expensive objective functions. In *Advances in Neural Information Processing Systems*,  
698 volume 28, 2015. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
699 2015/file/57c0531e13f40b91b3b0f1a30b529ald-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/57c0531e13f40b91b3b0f1a30b529ald-Paper.pdf).
- 700  
701 Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine  
learning algorithms. In *Advances in Neural Information Processing Systems*, volume 25,  
pp. 2951–2959, 2012. URL [https://proceedings.neurips.cc/paper\\_files/  
paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf).
- Arno Solin and Simo Särkkä. Hilbert space methods for reduced-rank Gaussian process regression.  
*Statistics and Computing*, 30(2):419–446, 2020. doi: 10.1007/s11222-019-09886-w.

- 702 Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias Seeger. Gaussian process  
703 optimization in the bandit setting: No regret and experimental design. In *Proceedings of the*  
704 *27th International Conference on Machine Learning*, volume 13, pp. 1015–1022, 2010. URL  
705 <https://icml.cc/Conferences/2010/papers/422.pdf>.  
706
- 707 S Surjanovic and D Bingham. Virtual library of simulation experiments: Test functions and datasets,  
708 2013. URL <http://www.sfu.ca/~ssurjano/optimization.html>.
- 709 Lloyd N. Trefethen. *Approximation Theory and Approximation Practice, Extended Edition*. So-  
710 ciety for Industrial and Applied Mathematics, Philadelphia, PA, 2019. ISBN 978-1-61197-  
711 593-2. doi: 10.1137/1.9781611975949. URL [https://people.maths.ox.ac.uk/  
712 trefethen/ATAP/](https://people.maths.ox.ac.uk/trefethen/ATAP/).
- 713 Ben Tu, Axel Gandy, Nikolas Kantas, and Behrang Shafei. Joint entropy  
714 search for multi-objective bayesian optimization. In *Advances in Neural In-*  
715 *formation Processing Systems*, volume 35, pp. 9922–9938, 2022. URL  
716 [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/  
717 4086fe59dc3584708468fba0e459f6a7-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/4086fe59dc3584708468fba0e459f6a7-Paper-Conference.pdf).
- 718
- 719 Jialei Wang, Scott C. Clark, Eric Liu, and Peter I. Frazier. Parallel Bayesian global optimization  
720 of expensive functions. *Operations Research*, 68(6):1850–1865, 2020. doi: 10.1287/opre.2019.  
721 1966. URL <https://doi.org/10.1287/opre.2019.1966>.
- 722
- 723 Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient Bayesian optimization. In  
724 *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 3627–  
725 3635, 2017. URL <https://proceedings.mlr.press/v70/wang17e.html>.
- 726
- 727 James Wilson, Frank Hutter, and Marc Deisenroth. Maximizing acquisition functions for  
728 Bayesian optimization. In *Advances in Neural Information Processing Systems*, volume 31, pp.  
729 9884–9895, 2018. URL [https://proceedings.neurips.cc/paper/2018/hash/  
730 498f2c21688f6451d9f5fd09d53edda7-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/498f2c21688f6451d9f5fd09d53edda7-Abstract.html).
- 731
- 732 James Wilson, Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc Deisenroth.  
733 Efficiently sampling functions from Gaussian process posteriors. In *Proceedings of the 37th*  
*International Conference on Machine Learning*, volume 119, pp. 10292–10302, 13–18 Jul 2020.  
734 URL <https://proceedings.mlr.press/v119/wilson20a.html>.
- 735
- 736 James T. Wilson, Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc Peter  
737 Deisenroth. Pathwise conditioning of gaussian processes. *Journal of Machine Learning Research*,  
738 22(105):1–47, 2021. URL <http://jmlr.org/papers/v22/20-1260.html>.
- 739
- 740 Huaiyu Zhu, C. K. I Williams, R Rohwer, and M Morciniec. Gaussian regression and optimal finite  
741 dimensional linear models. In *Neural Networks and Machine Learning*, 1998. URL <https://publications.aston.ac.uk/id/eprint/38366/>.  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A CHARACTERIZING THE LOCAL MINIMA OF A SEPARABLE FUNCTION

### A.1 PROOF OF PROPOSITION 1: A REPRESENTATION OF THE SET OF LOCAL MINIMA

Proposition 1 broadly applies to separable functions on a hypercube. Consider a separable function  $f(\mathbf{x}) = \prod_{i=1}^d f_i(x_i)$  with domain  $\mathcal{X} = \prod_{i=1}^d [\underline{x}_i, \bar{x}_i]$ , where  $f_i \in C^1([\underline{x}_i, \bar{x}_i]; \mathbb{R})$ . To simplify the discussion, we further assume that  $f_i$  is twice differentiable at its interior critical points  $\hat{x}_i$ . The gradient of  $f$  can be written as:

$$\nabla f(\mathbf{x}) = \left( f'_i(x_i) \cdot \prod_{j \neq i} f_j(x_j) \right)_{i=1}^d = \left( f(\mathbf{x}) \cdot \frac{f'_i(x_i)}{f_i(x_i)} \right)_{i=1}^d = f(\mathbf{x}) \cdot \mathbf{v}(\mathbf{x}), \quad (6)$$

where  $\mathbf{v}(\mathbf{x}) = (f'_i/f_i)_{i=1}^d = \left( \frac{d}{dx_i} \log f_i \right)_{i=1}^d$ . The Hessian of  $f$  can be written as:

$$\nabla^2 f(\mathbf{x}) = \text{diag} \left\{ f''_i(x_i) \prod_{j \neq i} f_j(x_j) \right\}_{i=1}^d + \left[ f'_i(x_i) f'_j(x_j) \prod_{k \neq i, j} f_k(x_k) \right]_{i \in d}^{j \neq i} = f(\mathbf{x}) \text{diag}(\mathbf{s} + \mathbf{v}\mathbf{v}^\top), \quad (7)$$

where  $\mathbf{s}(\mathbf{x}) = (f''_i/f_i - (f'_i/f_i)^2)_{i=1}^d = \left( \frac{d^2}{dx_i^2} \log f_i \right)_{i=1}^d$ .

An interior point  $\mathbf{x} \in \text{int } \mathcal{X} := \prod_{i=1}^d (\underline{x}_i, \bar{x}_i)$  is a strong local minimum of  $f$  if and only if  $\nabla f(\mathbf{x}) = 0$  and  $\nabla^2 f(\mathbf{x}) > 0$ . From eq. (6), the first condition is satisfied in any of the following three cases: (1)  $f_i(x_i) \neq 0$  and  $f'_i(x_i) = 0$  for all  $i \in \{1, \dots, d\}$ ; (2)  $f_i(x_i) = 0$  for exactly one  $i \in \{1, \dots, d\}$  and  $f'_i(x_i) = 0$ ; or (3)  $f_i(x_i) = 0$  for all  $i \in I \subseteq \{1, \dots, d\}$  where  $|I| \geq 2$ .

In case (1), the Hessian eq. (7) reduces to  $\nabla^2 f(\mathbf{x}) = f(\mathbf{x}) \cdot \text{diag}\{f''_i(x_i)/f_i(x_i)\}_{i=1}^d$ , which is positive definite if and only if one of the following holds: (i)  $f(\mathbf{x}) > 0$  and  $f_i(x_i)f''_i(x_i) > 0$ , for all  $i \in \{1, \dots, d\}$ ; or (ii)  $f(\mathbf{x}) < 0$  and  $f_i(x_i)f''_i(x_i) < 0$ , for all  $i \in \{1, \dots, d\}$ .

In case (2), the Hessian reduces to an all-zero matrix except for the  $i$ th diagonal entry:  $[\nabla^2 f(\mathbf{x})]_{i,i} = f''_i(x_i) \prod_{j \neq i} f_j(x_j)$ . Even if this entry is positive, the Hessian is still positive semi-definite, which means that there is a continuum of weak local minima:  $\{x_i\} \times \prod_{j \neq i} [\underline{x}_j, \bar{x}_j]$ . Besides, this case requires  $f_i$  and  $f'_i$  to have an identical root, which an event with probability zero.

In case (3), let  $g_i(r_i) := f_i(x_i + r_i)$  be a shifted version of  $f_i$ ,  $i \in \{1, \dots, d\}$ . Taylor expansion at  $\mathbf{r} = 0$  gives  $g_i(r_i) = 0 + g'_i(0)r_i + o(r_i)$  for all  $i \in I$  and  $g_j(r_j) = g_j(0) + O(r_j)$  for all  $j \notin I$ . We have  $g(\mathbf{r}) := \prod_{i=1}^d g_i(r_i) = c \prod_{i \in I} r_i + o(\prod_{i \in I} r_i) \cdot O(\prod_{j \notin I} r_j)$ , where  $c = \prod_{i \in I} g'_i(0) \cdot \prod_{j \notin I} g_j(0) \neq 0$ . This means that there is a continuum of saddle points:  $\{x_i\}_{i \in I} \times \prod_{j \notin I} [\underline{x}_j, \bar{x}_j]$ .

For a boundary point  $\mathbf{x} \in \partial \mathcal{X} := \mathcal{X} \setminus \text{int } \mathcal{X}$ , we partition the index set  $\{1, \dots, d\}$  into  $L, R$ , and  $I$  such that  $x_i = \underline{x}_i$  for all  $i \in L$ ,  $x_i = \bar{x}_i$  for all  $i \in R$ , and  $x_i \in (\underline{x}_i, \bar{x}_i)$  for all  $i \in I$ . Define  $\nabla_J := (\partial_j)_{j \in J}$  for any subset  $J$  of the indices. Then  $\mathbf{x}$  is a strong local minimum of  $f$  if and only if the following conditions hold: (a)  $\mathbf{x}$  is a strong local minimum in  $\{x_j\}_{j \notin I} \times \prod_{j \in I} [\underline{x}_j, \bar{x}_j]$ ; (b)  $\nabla_L f(\mathbf{x}) > 0$ ; and (c)  $\nabla_R f(\mathbf{x}) < 0$ .

Condition (a) holds if and only if  $\nabla_I f(\mathbf{x}) = 0$  and  $\nabla_I^2 f(\mathbf{x}) > 0$ . Based on the previous discussion on interior local minima, it is equivalent to: (i)  $f(\mathbf{x}) > 0$  and  $f_i(x_i)f''_i(x_i) > 0$ , for all  $i \in I$ ; or (ii)  $f(\mathbf{x}) < 0$  and  $f_i(x_i)f''_i(x_i) < 0$ , for all  $i \in I$ .

From eq. (6), condition (b) is equivalent to: (i)  $f(\mathbf{x}) > 0$  and  $f_i(x_i)f'_i(x_i) > 0$ , for all  $i \in L$ ; or (ii)  $f(\mathbf{x}) < 0$  and  $f_i(x_i)f'_i(x_i) < 0$ , for all  $i \in L$ .

Similarly, condition (c) is equivalent to: (i)  $f(\mathbf{x}) > 0$  and  $-f_i(x_i)f'_i(x_i) > 0$ , for all  $i \in R$ ; or (ii)  $f(\mathbf{x}) < 0$  and  $-f_i(x_i)f'_i(x_i) < 0$ , for all  $i \in R$ .

Summarizing the above discussions, we see that there is a unified way to identify the set  $\check{X}$  of all strong local minima of  $f$ , which is stated in Proposition 1. The discussion for the set  $\hat{X}$  of local maxima is the exactly the same, except that the signs are flipped. This also means that  $\hat{X}$  and  $\check{X}$  form a partition of the union  $\Xi^{(0)} \sqcup \Xi^{(1)}$  of the two tensor grids.

If  $f_i$  is not twice differentiable at some interior critical point  $x_i$ , we may replace  $f''_i(x_i) > 0$  with the statement that  $x_i$  is a strong local minimum of  $f_i$ , and replace  $f''_i(x_i) < 0$  with the statement

that  $x_i$  is a strong local maximum of  $f_i$ . The rest of the discussion still follows. In practice, the differentiability of the prior sample is not an issue, because it is almost always approximated by a finite sum of analytic functions, which is again analytic.

## A.2 NUMBER OF LOCAL MINIMA OF A SEPARABLE FUNCTION

In proposition 1, each set of candidate coordinates  $\Xi_i$  is partitioned into mixed type and mono type:

$$\Xi_i^{(1)} = \{\xi_{i,j} \in \Xi_i : f_i(\xi_{i,j})h_i(\xi_{i,j}) < 0\}, \quad \Xi_i^{(0)} = \{\xi_{i,j} \in \Xi_i : f_i(\xi_{i,j})h_i(\xi_{i,j}) > 0\}.$$

Another partition of  $\Xi_i$  is by the sign of the corresponding component function value:

$$\Xi_i^- = \{\xi_{i,j} \in \Xi_i : f_i(\xi_{i,j}) < 0\}, \quad \Xi_i^+ = \{\xi_{i,j} \in \Xi_i : f_i(\xi_{i,j}) > 0\}.$$

These two partitions create a finer partition of  $\Xi_i$  into four subsets:

$$\Xi_i^{-(1)} = \Xi_i^- \cap \Xi_i^{(1)}, \quad \Xi_i^{-(0)} = \Xi_i^- \cap \Xi_i^{(0)}, \quad \Xi_i^{+(1)} = \Xi_i^+ \cap \Xi_i^{(1)}, \quad \Xi_i^{+(0)} = \Xi_i^+ \cap \Xi_i^{(0)}.$$

Denote the sizes of mixed and mono type candidate coordinates as  $n_i^{(1)} = |\Xi_i^{(1)}|$  and  $n_i^{(0)} = |\Xi_i^{(0)}|$ , then the sizes of the two tensor grids  $\Xi^{(1)}$  and  $\Xi^{(0)}$  can be written as:

$$N^{(1)} := |\Xi^{(1)}| = \prod_{i=1}^d n_i^{(1)}, \quad N^{(0)} := |\Xi^{(0)}| = \prod_{i=1}^d n_i^{(0)}.$$

Define signed sums as the sums of signs of function values on the two tensor grids:

$$S^{(1)} := \sum_{\xi \in \Xi^{(1)}} \text{sign}(f(\xi)), \quad S^{(0)} := \sum_{\xi \in \Xi^{(0)}} \text{sign}(f(\xi)).$$

We now derive efficient formulas to calculate these signed sums, using  $S^{(1)}$  as an example. Denote each coordinate in  $\Xi_i^{(1)}$  as  $\xi_{i,j}^{(1)}$ . Denote each point in  $\Xi^{(1)}$  as  $\xi_J^{(1)} = (\xi_{i,J_i}^{(1)})_{i=1}^d$ , where multi-index  $J = (J_i)_{i=1}^d \in \Pi^{(1)} := \prod_{i=1}^d \{1, \dots, n_i^{(1)}\}$ . The signed sum  $S^{(1)}$  can be written as:

$$\begin{aligned} S^{(1)} &= \sum_{J \in \Pi^{(1)}} \text{sign}(f(\xi_J^{(1)})) = \sum_{J \in \Pi^{(1)}} \text{sign}\left(\prod_{i=1}^d f_i(\xi_{i,J_i}^{(1)})\right) \\ &= \sum_{J \in \Pi^{(1)}} \prod_{i=1}^d \text{sign}(f_i(\xi_{i,J_i}^{(1)})) = \prod_{i=1}^d \sum_{j=1}^{n_i^{(1)}} \text{sign}(f_i(\xi_{i,j}^{(1)})) \\ &= \prod_{i=1}^d \left[ \sum_{j=1}^{n_i^{(1)}} 1(f_i(\xi_{i,j}^{(1)}) > 0) - \sum_{j=1}^{n_i^{(1)}} 1(f_i(\xi_{i,j}^{(1)}) < 0) \right] = \prod_{i=1}^d \left[ |\Xi_i^{+(1)}| - |\Xi_i^{-(1)}| \right]. \end{aligned}$$

A formula for  $S^{(0)}$  can be derived analogously. Denote set sizes:

$$n_i^{-(1)} = |\Xi_i^{-(1)}|, \quad n_i^{-(0)} = |\Xi_i^{-(0)}|, \quad n_i^{+(1)} = |\Xi_i^{+(1)}|, \quad n_i^{+(0)} = |\Xi_i^{+(0)}|,$$

then the signed sums can be calculated as:

$$S^{(1)} = \prod_{i=1}^d (n_i^{+(1)} - n_i^{-(1)}), \quad S^{(0)} = \prod_{i=1}^d (n_i^{+(0)} - n_i^{-(0)}).$$

The sizes of negative and positive strong local minima of a separable function can be written as:

$$\check{N}^- := |\check{X}^-| = \sum_{\xi \in \Xi^{(1)}} 1(f(\xi) < 0) = \frac{1}{2}(N^{(1)} - S^{(1)}), \quad (8)$$

$$\check{N}^+ := |\check{X}^+| = \sum_{\xi \in \Xi^{(0)}} 1(f(\xi) > 0) = \frac{1}{2}(N^{(0)} + S^{(0)}).$$

Therefore, the size of the strong local minima of a separable function can be written as:

$$\check{N} := |\check{X}| = |\check{X}^-| + |\check{X}^+| = \frac{1}{2}(N^{(1)} + N^{(0)} - S^{(1)} + S^{(0)}). \quad (9)$$



## B ORDERING THE LOCAL MINIMA OF A SEPARABLE FUNCTION

### B.1 FILTERING A TENSOR GRID FOR HIGH ABSOLUTE VALUES OF A SEPARABLE FUNCTION

The step one in Section 3.4 is equivalent to the following problem: given coordinates  $Z_i = \{\zeta_{i,1}, \dots, \zeta_{i,t_i}\}$  and components values  $F_i = \{f_{i,1}, \dots, f_{i,t_i}\}$ ,  $i \in \{1, \dots, d\}$ , of a separable function  $f(\mathbf{x}) = \prod_{i=1}^d f_i(x_i)$ , find points  $\zeta$  such that  $|f(\zeta)|$  are the  $k$  largest in the tensor grid  $Z = \prod_{i=1}^d Z_i$ .

Because  $\log |f(\mathbf{x})| = \log |\prod_{i=1}^d f_i(x_i)| = \sum_{i=1}^d \log |f_i(x_i)|$ , we can solve this problem as follows: define two-dimensional arrays  $F = [F_1, \dots, F_d]$  and  $A = \log |F|$ , solve  $S = \text{maxk\_sum}(A, k)$ , and return  $\{\zeta = (\zeta_{1,I_1}, \dots, \zeta_{d,I_d}) : I \in S\}$ . Here the `maxk_sum` algorithm finds the combinations from  $A$  that gives the  $k$  largest sums, which is described next.

### B.2 TOP COMBINATIONS WITH THE LARGEST SUMS

Consider this problem: given a two-dimensional array  $A = [\mathbf{a}_1, \dots, \mathbf{a}_d]$ ,  $\mathbf{a}_i = [a_{i,1}, \dots, a_{i,t_i}]$ , with  $a_{i,1} \geq \dots \geq a_{i,t_i}$ ,  $i \in \{1, \dots, d\}$ , find  $k$  multi-indices of the form  $I = [I_1, \dots, I_d]$  such that the sums  $s_I := \sum_{i=1}^d a_{i,I_i}$  are the  $k$  largest among all combinations  $I \in \prod_{i=1}^d \{1, \dots, t_i\}$ .

An exhaustive search is intractable because the number of all possible combinations grows exponentially as  $\prod_{i=1}^d t_i$ . Instead, we use a min-heap to efficiently keep track of the top  $k$  combinations. A min-heap is a complete binary tree, where each node is no greater than its children. The operations of inserting an element and removing the smallest element from a min-heap can be done in logarithmic time. Algorithm 2 gives a procedure to solve the above problem using min-heaps.

---

**Algorithm 2** `maxk_sum`: Combinations with the  $k$  largest sums

---

**Input:** two-dimensional array  $A$ ; number of top combinations  $k$ .

- 1: Make the array nonpositive by replacing  $\mathbf{a}_i$  with  $\mathbf{a}_i - a_{i,1}\mathbf{1}$  for  $i = 1, \dots, d$ .
- 2: Create a min-heap by adding the elements of  $\mathbf{a}_1$ , each considered a combination of length one: index  $I_1$ , key  $a_{1,I_1}$ .
- 3: At stage  $i = 2, \dots, d$ : create a new min-heap consisting of length- $i$  combinations by adding each element in  $\mathbf{a}_i$  to each combination in the min-heap at the previous stage: index  $[I_1, \dots, I_i]$ , key  $\sum_{j=1}^i a_{j,I_j}$ . The size of the min-heap at each stage is capped at  $k$  by popping the smallest sum from the min-heap when necessary.

**Output:** combinations in the min-heap at stage  $d$ .

---

This algorithm has time complexity  $\mathcal{O}(tk \log k)$ , where  $t = \sum_{i=1}^d t_i \ll \prod_{i=1}^d t_i$ , and space complexity  $\mathcal{O}(dk)$ . In TS-roots, the cost of `maxk_sum` is small compared with the gradient-based multistart optimization of the posterior sample.

## C ALGORITHMS FOR TS-ROOTS

### C.1 SPECTRAL SAMPLING OF SEPARABLE GAUSSIAN PROCESS PRIORS

Per Mercer’s theorem on probability spaces (see e.g., Rasmussen & Williams (2006), Sec 4.3), any positive definite covariance function that is essentially bounded with respect to some probability measure  $\mu$  on a compact domain  $\mathcal{X}$  has a spectral representation  $\kappa(\mathbf{x}, \mathbf{x}') = \sum_{k=0}^{\infty} \lambda_k \phi_k(\mathbf{x}) \phi_k(\mathbf{x}')$ , where  $(\lambda_k, \phi_k(\mathbf{x}))$  is a pair of eigenvalue and eigenfunction of the kernel integral operator. The corresponding GP prior can be written as  $f_\omega(\mathbf{x}) = \sum_{k=0}^{\infty} w_k \sqrt{\lambda_k} \phi_k(\mathbf{x})$ , where  $w_k \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$  are independent standard Gaussian random variables. Similar spectral representations exist per Bochner’s theorem, which may have efficient discretizations (Solin & Särkkä, 2020; Mutny & Krause, 2018).

Given spectral representations of the univariate component functions of a separable Gaussian Process prior, we can accurately approximate the prior sample as:

$$f_{\omega}(\mathbf{x}) = \prod_{i=1}^d f_i(x_i; \omega_i), \quad f_i(x_i; \omega_i) \approx \sum_{k=0}^{N_i-1} w_{i,k} \sqrt{\lambda_{i,k}} \phi_{i,k}(x_i). \quad (10)$$

Here  $N_i$  is selected for each variate such that  $\lambda_{i,N_i-1}/\lambda_{i,1} \leq \eta_i$ , where  $\eta_i$  is sufficiently small (see Appendix H for the value used in this study). Using spectral representations of the univariate components as in eq. (10) is much more efficient than directly using a spectral representation of the separable GP prior, because the former uses  $\sum_{i=1}^d N_i$  univariate terms to exactly represent  $\prod_{i=1}^d N_i$  multivariate terms in the latter.

**Spectrum of the Squared Exponential Covariance Function.** The univariate squared exponential (SE) covariance function can be written as  $\kappa(x, x'; l) = \exp(-\frac{1}{2}s^2)$ , where the relative distance  $s = |x - x'|/l$  and length scale  $l \in (0, \infty)$ . The spectral representation of such a covariance function per Mercer’s theorem is  $\kappa(x, x') = \sum_{k=0}^{\infty} \lambda_k \phi_k(x) \phi_k(x')$ . With a Gaussian measure  $\mu = \mathcal{N}(0, \sigma^2)$  over the domain  $\mathcal{X} = \mathbb{R}$ , we can write the eigenvalues  $\lambda_k$  and eigenfunctions  $\phi_k(x)$  of the kernel integral operator as follows. (See e.g., Zhu et al. (1998) Sec. 4 and Gradshteyn & Ryzhik (2014) 7.374 eq. 8.)

Define constants  $a = (2\sigma^2)^{-1}$ ,  $b = (2l)^{-1}$ ,  $c = \sqrt{a^2 + 4ab}$ , and  $A = \frac{1}{2}a + b + \frac{1}{2}c$ . For  $k \in \mathbb{N}$ , the  $k$ th eigenvalue is  $\lambda_k = \sqrt{\frac{a}{A}} \left(\frac{b}{A}\right)^k$  and the corresponding eigenfunction is  $\phi_k(x) = \left(\frac{\pi c}{a}\right)^{1/4} \psi_k(\sqrt{c}x) \exp\left(\frac{1}{2}ax^2\right)$ , where  $\psi_k(x) = (\pi^{1/2} 2^k k!)^{-1/2} H_k(x) \exp(-\frac{1}{2}x^2)$  and  $H_k(x)$  the  $k$ th-order Hermite polynomial defined by  $H_k(x) = (-1)^k \exp(x^2) \frac{d^k}{dx^k} \exp(-x^2)$ .

Figure 6 shows approximations to the SE covariance function by truncated spectral representations with the first  $N$  eigenpairs and by random Fourier features (Rahimi & Recht, 2007) with  $N$  basis functions. The spectral representation per Mercer’s theorem converges quickly to the true covariance function, while the random Fourier features representation requires a large number of basis functions and is inaccurate for  $N < 1000$ .

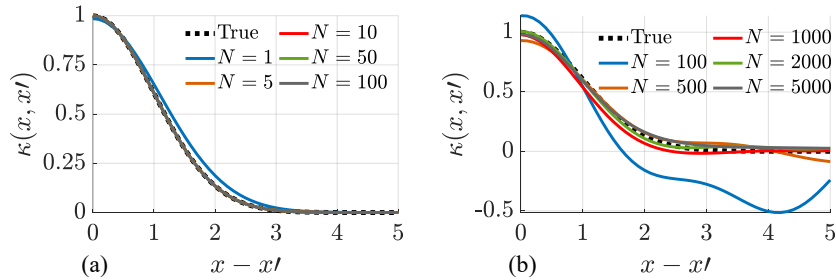


Figure 6: Approximate SE covariance functions from (a) the spectral representation per Mercer’s theorem with the first  $N$  eigenpairs and (b) the random Fourier features representation with  $N$  basis functions. The plots are generated for  $l = 1$ .

## C.2 UNIVARIATE GLOBAL ROOTFINDING

Algorithm 3 outlines a method for univariate global rootfinding on an interval by solving an eigenvalue problem. When the orthogonal polynomial basis is the Chebyshev polynomials, the corresponding comrade matrix is called a colleague matrix, and we have the following theorem:

**Theorem 1** Let  $p(x) = \sum_{k=0}^m a_k T_k(x)$ ,  $a_m \neq 0$ , be a polynomial of degree  $m$ , where  $T_k$  is the  $k$ th Chebyshev polynomial and  $a_k$  is the corresponding weight. The roots of  $p(x)$  are the eigenvalues of



**Algorithm 4** minsort: Best local minima of a separable function

---

1026 **Algorithm 4** minsort: Best local minima of a separable function

1027

1028 **Input:** separable function  $f(\mathbf{x}) = \prod_{i=1}^d f_i(x_i)$ ; set size  $n_o$ ; buffer coefficient  $\alpha$  (defaults to 3).

1029 1:  $f_i(x_i) \leftarrow \text{chebfun}(f_i(x_i))$ ,  $i = 1, \dots, d$   $\triangleright$  Construct chebfuns for univariate components

1030  $f_i'(x_i) \leftarrow \text{diff}(f_i(x_i))$ ;  $f_i''(x_i) \leftarrow \text{diff}(f_i'(x_i))$   $\triangleright$  Compute first and second derivatives

1031 2:  $\{\xi_{i,j}\}_{j=1}^{r_i} \leftarrow \text{roots}(f_i'(x_i))$ ,  $i = 1, \dots, d$   $\triangleright$  Univariate global rootfinding

1032  $\{\xi_{i,0}\} \leftarrow \underline{x}_i$ ;  $\{\xi_{i,r_i+1}\} \leftarrow \bar{x}_i$   $\triangleright$  Include interval lower and upper bounds

1033  $\xi_i \leftarrow [\xi_{i,0}, \xi_{i,1}, \dots, \xi_{i,r_i}, \xi_{i,r_i+1}]^\top$   $\triangleright$  Candidate coordinate values  $\{\Xi_i\}$

1034 3:  $\mathbf{f}_i \leftarrow f_i(\xi_i)$ ,  $i = 1, \dots, d$   $\triangleright$  Univariate function values

1035  $h_{i,j} \leftarrow f_i''(\xi_{i,j})$ ,  $j = 1, \dots, r_i$   $\triangleright$  Univariate second derivatives at critical points

1036  $h_{i,0} \leftarrow f_i''(\xi_{i,0})$ ;  $h_{i,r_i+1} \leftarrow -f_i''(\xi_{i,r_i+1})$   $\triangleright$  Univariate inward derivatives at interval ends

1037 4:  $J_i \leftarrow (\mathbf{f}_i \circ \mathbf{h}_i > 0)$ ;  $P_i \leftarrow (\mathbf{f}_i > 0)$   $\triangleright$  Boolean vectors of sign parity and positivity

1038 5:  $\xi_i^{(0)} \leftarrow \xi_i(J_i)$ ;  $\xi_i^{(1)} \leftarrow \xi_i(\neg J_i)$   $\triangleright$  Mono and mixed type candidate coordinates:  $\Xi_i^{(0)}, \Xi_i^{(1)}$

1039  $\mathbf{f}_i^{(0)} \leftarrow \mathbf{f}_i(J_i)$ ;  $\mathbf{f}_i^{(1)} \leftarrow \mathbf{f}_i(\neg J_i)$   $\triangleright$  Values at mono and mixed type candidate coordinates

1040 6:  $n_i^{(0)} \leftarrow \text{sum}(J_i)$ ;  $n_i^{(1)} \leftarrow \text{sum}(\neg J_i)$

1041  $n_i^{+(0)} \leftarrow \text{sum}(P_i \& J_i)$ ;  $n_i^{-(0)} \leftarrow \text{sum}((\neg P_i) \& J_i)$

1042  $n_i^{+(1)} \leftarrow \text{sum}(P_i \& (\neg J_i))$ ;  $n_i^{-(1)} \leftarrow \text{sum}((\neg P_i) \& (\neg J_i))$

1043  $N^{(0)} \leftarrow \prod_{i=1}^d n_i^{(0)}$ ;  $N^{(1)} \leftarrow \prod_{i=1}^d n_i^{(1)}$   $\triangleright$  Sizes of tensor grids

1044  $S^{(0)} \leftarrow \prod_{i=1}^d (n_i^{+(0)} - n_i^{-(0)})$ ;  $S^{(1)} \leftarrow \prod_{i=1}^d (n_i^{+(1)} - n_i^{-(1)})$   $\triangleright$  Signed sums

1045 7: **if**  $n_o \leq \check{N}^- = \frac{1}{2}(N^{(1)} - S^{(1)})$  **then**

1046

1047 8:  $[\mathbf{s}, \mathbf{I}] \leftarrow \text{maxk\_sum}(\{\log(|\mathbf{f}_i^{(1)}|)\}_{i=1}^d, \alpha n_o)$   $\triangleright$  The  $\alpha n_o$  largest  $|f|$  in  $\Xi^{(1)}$

1048

1049 9:  $\mathbf{I} \leftarrow \mathbf{I}[f(\xi^{(1)}(\mathbf{I})) < 0, :]$   $\triangleright$  Multi-indices of best negative local minima

1050 10:  $[\mathbf{b}, I] \leftarrow \text{mink}(f(\xi^{(1)}(\mathbf{I})), n_o)$   $\triangleright$  The  $n_o$  smallest  $f$  in  $\check{X}^-$

1051 11:  $S_o \leftarrow S_o^- = \xi^{(1)}(\mathbf{I}[I, :])$

1052 12: **else**

1053 13:  $\Pi^{(1)} \leftarrow \prod_{i=1}^d \{1, \dots, n_i^{(1)}\}$   $\triangleright$  Matrix of index combinations

1054 14:  $\check{\mathbf{I}}^- \leftarrow \Pi^{(1)}[f(\xi^{(1)}(\Pi^{(1)})) < 0, :]$   $\triangleright$  Multi-indices of negative local minima

1055 15:  $[\mathbf{b}, I] \leftarrow \text{sort}(f(\xi^{(1)}(\check{\mathbf{I}}^-)))$   $\triangleright$  Sort values in ascending order

1056 16:  $\check{X}^- \leftarrow \xi^{(1)}(\check{\mathbf{I}}^-[I, :])$   $\triangleright$  Negative local minima

1057 17: **if**  $n_o \leq \check{N} = \frac{1}{2}(N^{(1)} - S^{(1)} + N^{(0)} + S^{(0)})$  **then**

1058 18:  $[\mathbf{s}, \mathbf{I}] \leftarrow \text{maxk\_sum}(\{\log(|\mathbf{f}_i^{(0)}|)\}_{i=1}^d, \alpha(n_o - \check{N}^-))$   $\triangleright$  Largest  $|f|$  in  $\Xi^{(0)}$

1059

1060 19:  $\mathbf{I} \leftarrow \mathbf{I}[f(\xi^{(0)}(\mathbf{I})) > 0, :]$   $\triangleright$  Multi-indices of best positive local minima

1061 20:  $[\mathbf{b}, I] \leftarrow \text{mink}(f(\xi^{(0)}(\mathbf{I})), n_o - \check{N}^-)$   $\triangleright$  The  $n_o - \check{N}^-$  smallest  $f$  in  $\check{X}^+$

1062 21:  $S_o \leftarrow \check{X}^- \cup S_o^+$ ,  $S_o^+ = \xi^{(0)}(\mathbf{I}[I, :])$

1063 22: **else**

1064 23:  $\Pi^{(0)} \leftarrow \prod_{i=1}^d \{1, \dots, n_i^{(0)}\}$   $\triangleright$  Matrix of index combinations

1065 24:  $\check{\mathbf{I}}^+ \leftarrow \Pi^{(0)}[f(\xi^{(0)}(\Pi^{(0)})) > 0, :]$   $\triangleright$  Multi-indices of positive local minima

1066 25:  $[\mathbf{b}, I] \leftarrow \text{sort}(f(\xi^{(0)}(\check{\mathbf{I}}^+)))$   $\triangleright$  Sort values in ascending order

1067 26:  $S_o \leftarrow \check{X}^- \cup \check{X}^+$ ,  $\check{X}^+ = \xi^{(0)}(\check{\mathbf{I}}^+[I, :])$   $\triangleright$  All local minima

1068 27: **end if**

1069 28: **end if**

1070 **Output:**  $S_o$   $\triangleright$  Candidate exploration set: smallest  $n_o$  local minima in ascending order

---

1074 representation (Appendix C.1); and (ii) the canonical basis  $\kappa_{\cdot, n}(\cdot)$ , which costs  $\mathcal{O}(dn)$  flops. When

1075 the data size  $n$  is large, we can pre-filter the observed locations  $X$  by the observations  $\mathbf{y}$ , which is a

1076 good estimate of  $\tilde{f}(X)$  depending on the observation noise. We assume that the number of observed

1077 locations after filtering is at most comparable to  $n_e$ . The cost of task (2) is thus  $\mathcal{O}(n_o dn)$  flops.

1078 Now consider task (3). Evaluating the gradient of  $\tilde{f}(\cdot)$  involves evaluating the gradients of  $f(\cdot)$  and

1079  $\kappa_{\cdot, n}(\cdot)$ . Since both  $f(\cdot)$  and  $\kappa_{\cdot, n}(\cdot)$  are separable functions, their gradients can be computed at a

**Algorithm 5** Decoupled sampling of Gaussian process posterior

**Input:** eigenpairs  $\{(\lambda_{i,k}, \phi_{i,k}(x))\}_{i=1, \dots, d}^{k=0, \dots, N_i-1}$ , data  $\mathcal{D} = \{(\mathbf{x}^j, y^j)\}_{j=1}^n$ , covariance matrix  $\mathbf{C} = \mathbf{K}_{n,n} + \mathbf{\Sigma}$ , canonical basis  $\boldsymbol{\kappa}_{\cdot, n}(\mathbf{x}) = (\kappa(\mathbf{x}, \mathbf{x}^j))_{j=1}^n$ .

- 1:  $w_{i,k} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$  ▷ Random coefficients for the prior sample
- 2:  $f_\omega(\mathbf{x}) = \prod_{i=1}^d \sum_{k=0}^{N_i-1} w_{i,k} \sqrt{\lambda_{i,k}} \phi_{i,k}(x_i)$  ▷ Approximate prior sample
- 3:  $\mathbf{f}_n \leftarrow [f_\omega(\mathbf{x}^1), \dots, f_\omega(\mathbf{x}^n)]^\top$  ▷ Values of prior sample at observed locations
- 4:  $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \mathbf{\Sigma})$  ▷ Random noise for the posterior sample
- 5:  $\mathbf{v} \leftarrow \mathbf{C}^{-1}(\mathbf{y} - \mathbf{f}_n - \boldsymbol{\varepsilon})$  ▷ Linear solve via factorization (e.g., Cholesky or SVD)

**Output:**  $\tilde{f}_\omega(\mathbf{x}) = f_\omega(\mathbf{x}) + \mathbf{v}^\top \boldsymbol{\kappa}_{\cdot, n}(\mathbf{x})$  ▷ Approximate posterior sample

cost comparable to evaluating their function values. Let  $N_{\text{grad}}$  be the number of gradient evaluations required by the gradient-based optimizer. The cost of task (3) is thus  $\mathcal{O}((n_e + n_x)N_{\text{grad}}dn)$  flops.

For task (1), the cost of the `minsort` algorithm is dominated by: (i)  $d$  calls to `chebfun`, which evaluates the univariate components  $f_i(\cdot)$  at a number of points depending on their complexity; (ii)  $d$  calls to `roots` (Algorithm 3), which scales linearly with the polynomial degree  $m$  of the `chebfun` object, itself dependent on the complexity of  $f_i(\cdot)$ ; and (iii) at most one call to `maxk_sum` (Algorithm 2) at a cost of  $\mathcal{O}(tn_o \log n_o)$ , where  $t = \sum_{i=1}^d t_i$  and  $t_i$  is two plus the number of critical points of  $f_i(\cdot)$  which depends on the complexity of  $f_i(\cdot)$ . The complexity of  $f_i(\cdot)$ , for the SE kernel for example (see Appendix C.1), can be quantified as the inverse length scale  $\theta_i = 1/l_i$ . We may define an average complexity as  $\theta = \frac{1}{d} \sum_{i=1}^d \theta_i$ . The cost of task (1) is thus  $\mathcal{O}(d\theta n_o \log n_o)$  flops.

As explained in Appendix D, we can set  $n_e$  and  $n_x$  to small values and  $n_o$  to a moderate value, independent of  $\tilde{f}(\cdot)$  and thus independent of  $d$ ,  $n$  and  $\theta$ . The overall cost of `TS-ROOTS` thus scales as  $\mathcal{O}(dn + d\theta)$ , which is linear in the input dimension  $d$ .

## D MINIMUM SIZE OF EXPLORATION AND EXPLOITATION SETS

We conduct an empirical experiment to determine minimal values for  $n_e$  and  $n_x$  of `TS-roots` algorithm, which are the sizes of  $S_e$  and  $S_x$ , respectively. From this experiment, we also recommend a value for  $n_o$ , which is the size of  $S_o$ . Recall that points in  $S_o$  are sorted in ascending order of prior sample values, while those in  $S_e$  and  $S_x$  are sorted in ascending order of posterior sample values.

Let  $I_e$  and  $I_x$  be the sets of indices of points in  $S_e$  and  $S_x$  that converge to the best local minimum of the posterior sample in each optimization iteration, respectively. Let  $I_o$  be the set of indices of points in  $S_o$  associated with  $S_e$ . Our hypothesis is that we have a high chance of finding a small index value in either  $I_e$  or  $I_x$ . If this hypothesis is confirmed, then we can set both  $n_e$  and  $n_x$  at very small values, which significantly accelerate the inner-loop optimization. To confirm our hypothesis, we employ the following two steps. First, we set  $n_e$  and  $n_x$  at large values to mimic the effect of removing the set size limits, ensuring accurate solution of the global optimization problem. Then, we show that we have a high chance of finding a small index value from  $I_e$  and/or  $I_x$  in each optimization iteration.

We test our hypothesis on the 2D Schwefel, 4D Rosenbrock, 10D Levy, 16D Ackley, 16D Powell functions. We set  $n_o = 5000$ ,  $n_e = n_x = 1000$ , and  $\alpha = 3$  (buffer coefficient). The left and middle columns of Figure 7 show the smallest index values and the variation of index values from  $I_e$  and/or  $I_x$  of starting points that converge to the best local minimum  $\mathbf{x}^*$  of the posterior sample path in each optimization iteration. The left column also plots the index values from  $I_o$  corresponding to the smallest index values from  $I_e$ , if exists. The right column shows the histograms of the smallest index values from  $I_e$  and  $I_x$  for all iterations considered. These results show that we have a high chance of finding a small index value from  $I_e$  and/or  $I_x$  in each iteration. This confirms our hypothesis. In fact, using the first point in  $S_e$  and the first point in  $S_e$ —only two points—we can discover the global optimum most of the time. Interestingly, the smallest index values appear largely independent of both the optimization iteration and the input dimension. Furthermore, the results suggest that it is safe to set  $n_o = 500$ ; and for almost exact global optimization, we suggest setting  $n_e = 25$  and  $n_x = 50$ .

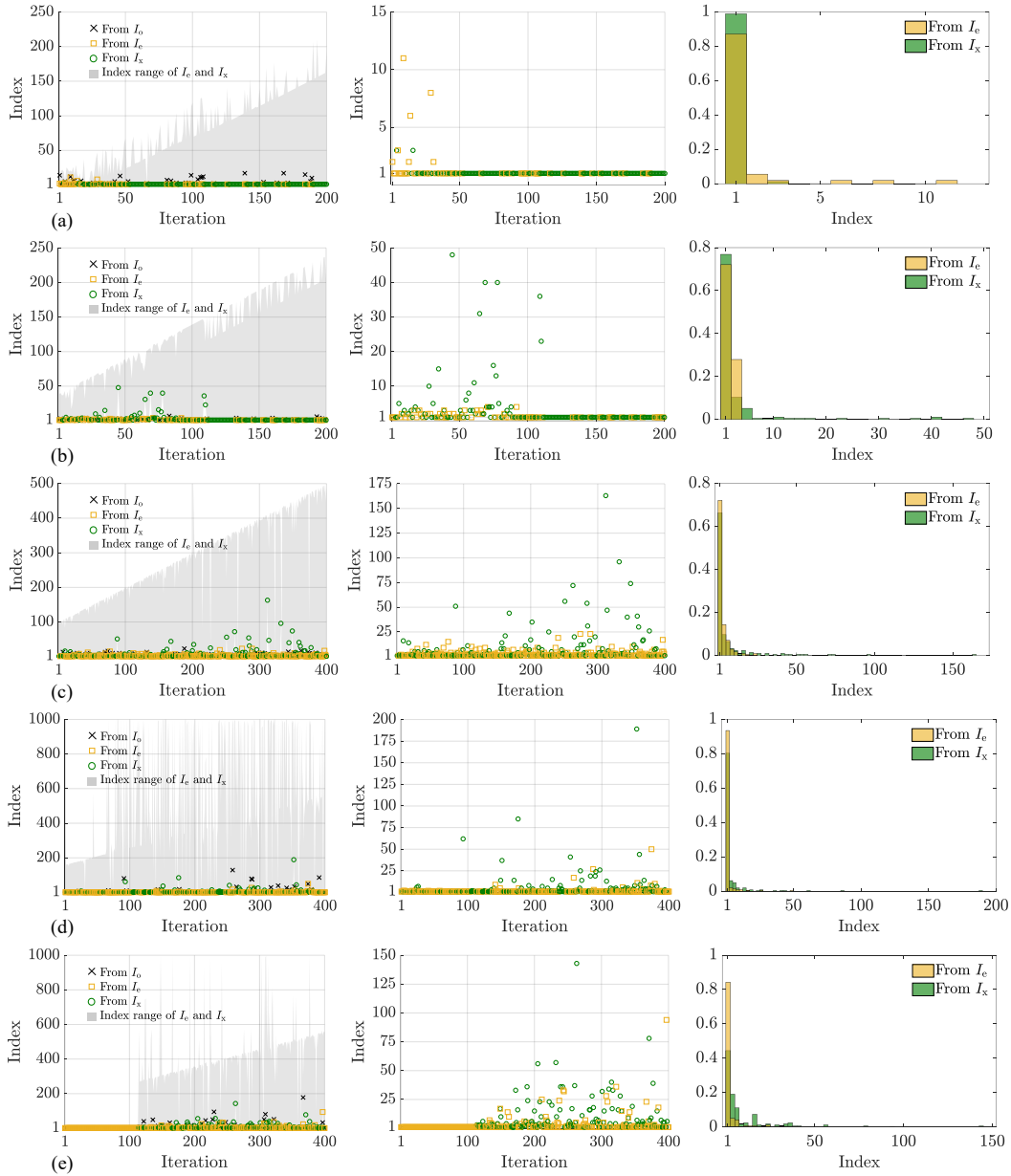


Figure 7: *Left column*: Minimum index values and index range of  $I_e$  and/or  $I_x$  for starting points that converge to the best local minimum  $\mathbf{x}^*$  of posterior sample in each optimization iteration, and index values of  $I_0$  associated with minimum index values from  $I_e$ . *Middle column*: Zoom-in plots of index values. *Right column*: Histogram of the minimum index values. (a) 2D Schwefel, (b) 4D Rosenbrock, (c) 10D Levy, (d) 16D Ackley, (e) 16D Powell functions.

## E BAYESIAN OPTIMIZATION VIA THOMPSON SAMPLING

A general procedure for sequential optimization is given in Algorithm 6. The initial dataset  $\mathcal{D}^0$  can either be empty or contain some observations. In the latter case we can write  $\mathcal{D}^0 = \{(\mathbf{x}^i, y^i)\}_{i=1}^{n_0}$ , where  $n_0 \in \mathbb{N}_{>0}$ . Three components of this algorithm can be customized: the observation model  $\text{Observe}(\mathbf{x})$ , the optimization policy  $\text{Policy}(\mathcal{D})$ , and the termination condition.

BO can be seen as an optimization policy for sequential optimization. A formal procedure is given in Algorithm 7. Three components of this algorithm can be customized: the prior probabilistic model

**Algorithm 6** Sequential optimization (Garnett, 2023)

---

**Input:** initial dataset  $\mathcal{D}^0$

- 1:  $k \leftarrow 1$
- 2: **repeat**
- 3:    $\mathbf{x}^k \leftarrow \text{Policy}(\mathcal{D}^{k-1})$
- 4:    $y^k \leftarrow \text{Observe}(\mathbf{x}^k)$
- 5:    $\mathcal{D}^k \leftarrow \mathcal{D}^{k-1} \cup \{(\mathbf{x}^k, y^k)\}$
- 6: **until** termination condition reached

**Output:**  $\mathcal{D}$

---

**Algorithm 7** Bayesian optimization policy

---

**Input:** a prior stochastic process  $f$  for the objective function  $f_{\text{true}}$ , current dataset  $\mathcal{D}^{k-1}$

- 1: determine the posterior  $f^k := f|\mathcal{D}^{k-1}$
- 2: derive an acquisition function  $\alpha^k(\mathbf{x})$  from  $f^k$
- 3: global optimization  $\mathbf{x}^k \leftarrow \arg \min_{\mathbf{x} \in \mathcal{X}} \alpha^k(\mathbf{x})$

**Output:**  $\mathbf{x}^k$

---

$f$ , the acquisition function  $\alpha$ , and the global optimization algorithm. Any probabilistic model of the objective function  $f_{\text{true}}$  can be seen as a probability distribution on a function space, and the prior  $f$  is usually specified as a stochastic process such as a GP. The acquisition function  $\alpha$  derived from the posterior  $f|\mathcal{D}$  can be either deterministic—such as EI and LCB—or stochastic, such as GP-TS. To simplify notation, we state the global optimization problem of  $\alpha(\mathbf{x})$  as minimization rather than maximization. The two problems are the same with a change of sign to the objective.

When applied to BO, GP-TS generates a random acquisition function simply by sampling the posterior model. That is, given the posterior  $f^k$  at the  $k$ th BO iteration, the GP-TS acquisition function is a random function:  $\alpha^k(\mathbf{x}) \sim f^k$ .

## F BENCHMARK FUNCTIONS

The analytical expressions for the benchmark functions used in Section 6 are given below. The global solutions of these functions are detailed in (Surjanovic & Bingham, 2013).

### Schwefel Function:

$$f(\mathbf{x}) = 418.9829d - \sum_{i=1}^d x_i \sin\left(\sqrt{|x_i|}\right). \quad (12)$$

This function is evaluated on  $\mathcal{X} = [-500, 500]^d$  and has a global minimum  $f^* := f(\mathbf{x}^*) = 0$  at  $\mathbf{x}^* = [420.9687, \dots, 420.9687]^\top$ . This function is  $C^1$  at  $\mathbf{x} = 0$ .

### Rosenbrock Function:

$$f(\mathbf{x}) = \sum_{i=1}^{d-1} [100(x_{i+1} - x_i^2)^2 + (x_i - 1)^2]. \quad (13)$$

This function is evaluated on  $\mathcal{X} = [-5, 10]^d$  and has a global minimum  $f^* = 0$  at  $\mathbf{x}^* = [1, \dots, 1]^\top$ .

### Levy Function:

$$f(\mathbf{x}) = \sin^2(\pi w_1) + \sum_{i=1}^{d-1} (w_i - 1)^2 [1 + 10 \sin^2(\pi w_i + 1)] + (w_d - 1)^2 [1 + \sin^2(2\pi w_d)], \quad (14)$$

where  $w_i = 1 + \frac{x_i - 1}{4}$ ,  $i = 1, \dots, d$ . This function is evaluated on  $\mathcal{X} = [-10, 10]^d$  and has a global minimum  $f^* = 0$  at  $\mathbf{x}^* = [1, \dots, 1]^\top$ .

**Ackley Function:**

$$f(\mathbf{x}) = -a \exp \left( -b \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2} \right) - \exp \left( \frac{1}{d} \sum_{i=1}^d \cos(cx_i) \right) + a + \exp(1), \quad (15)$$

where  $a = 20$ ,  $b = 0.2$ , and  $c = 2\pi$ . This function is evaluated on  $\mathcal{X} = [-10, 10]^d$  and has a global minimum  $f^* = 0$  at  $\mathbf{x}^* = [0, \dots, 0]^\top$ . This function not differentiable at  $\mathbf{x}^*$ .

**Powell Function:**

$$f(\mathbf{x}) = \sum_{i=1}^{d/4} \left[ (x_{4i-3} + 10x_{4i-2})^2 + 5(x_{4i-1} - x_{4i})^2 + (x_{4i-2} - 2x_{4i-1})^4 + 10(x_{4i-3} - x_{4i})^4 \right]. \quad (16)$$

This function is evaluated on  $\mathcal{X} = [-4, 5]^d$  and has a global minimum  $f^* = 0$  at  $\mathbf{x}^* = [0, \dots, 0]^\top$ .

**6d Hartmann Function:**

$$f(\mathbf{x}) = - \sum_{i=1}^4 a_i \exp \left( - \sum_{j=1}^6 A_{ij} (x_j - P_{ij})^2 \right), \quad (17)$$

where

$$\mathbf{a} = [1, 1.2, 3, 3.2]^\top, \quad (18a)$$

$$\mathbf{A} = \begin{bmatrix} 10 & 3 & 17 & 3.5 & 1.7 & 8 \\ 0.05 & 10 & 17 & 0.1 & 8 & 14 \\ 3 & 3.5 & 1.7 & 10 & 17 & 8 \\ 17 & 8 & 0.05 & 10 & 0.1 & 14 \end{bmatrix}, \quad (18b)$$

$$\mathbf{P} = 10^{-4} \begin{bmatrix} 1312 & 1696 & 5569 & 124 & 8283 & 5886 \\ 2329 & 4135 & 8307 & 3736 & 1004 & 9991 \\ 2348 & 1451 & 3522 & 2883 & 3047 & 6650 \\ 4047 & 8828 & 8732 & 5743 & 1091 & 381 \end{bmatrix} \quad (18c)$$

This function is evaluated on  $\mathcal{X} = [0, 1]^6$  and has a global minimum  $f^* = -3.32237$  at  $\mathbf{x}^* = [0.20169, 0.150011, 0.476874, 0.275332, 0.311625, 0.6573]^\top$ . The rescaled version  $\tilde{f}(\mathbf{x}) = \frac{f(\mathbf{x}) - 2.58}{1.94}$  (Picheny et al., 2013) is used in the experiments.

**G TEN-BAR TRUSS**

Consider a ten-bar truss shown in Figure 8. The truss has ten members and is subjected to vertical load  $P_1 = 60$  kN at node 2, vertical load  $P_2 = 40$  kN at node 3, and horizontal load  $P_3 = 40$  kN at node 3. The Young's modulus of the truss material  $E = 200$  GPa. The length parameter  $L = 1$  m. Let  $A(\mathbf{x}) = \sum_{i=1}^{10} x_i$  and  $\delta_3(\mathbf{x})$  denote the total area of the cross-sectional areas of the truss members and the vertical displacement at node 3, respectively, where  $\mathbf{x} = [x_1, \dots, x_{10}]^\top$  is the vector of cross-sectional areas of the truss members. The optimization problem formulated for the truss is to minimize both  $A(\mathbf{x})$  and  $\delta_3(\mathbf{x})$ . Since  $A(\mathbf{x})$  and  $\delta_3(\mathbf{x})$  are competing, we define the objective function as a weight-sum of  $A(\mathbf{x})$  and  $\delta_3(\mathbf{x})$ , such that

$$f(\mathbf{x}) = w_1 \frac{A(\mathbf{x})}{A_{\max}} + w_2 \frac{\delta_3(\mathbf{x})}{\delta_{\max}}, \quad (19)$$

where  $\mathbf{x} \in [1, 20]^{10}$  cm<sup>2</sup>,  $w_1 = 0.6$ ,  $w_2 = 0.4$ ,  $A_{\max} = 200$  cm<sup>2</sup>, and  $\delta_{\max} = 3$  cm.

**H EXPERIMENTAL DETAILS**

**Data Generation.** We generate 20 initial datasets for each problem. The input observations are randomly generated using the Latin hypercube sampling (Owen, 1992) within  $[-1, 1]^d$ , where  $d$  represents the number of input variables. The normalized input observations are transformed into their real spaces to evaluate the corresponding objective function values which are then standardized using the  $z$ -score for processing optimization. Each BO method in comparison starts from each of the generated datasets.



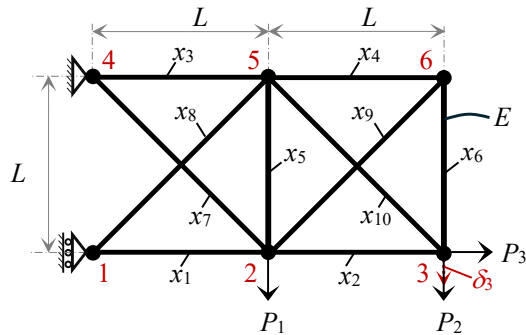


Figure 8: Ten-bar truss. Cross-sectional areas of ten truss members are the input variables  $x_i$ ,  $i \in \{1, \dots, 10\}$ . Known parameters include length  $L$ , Young’s modulus of truss material  $E$ , and external loads  $P_j$ ,  $j \in \{1, 2, 3\}$ . The vertical displacement at node 3 is denoted as  $\delta_3$ .

**Key Parameters for TS-roots and other BO Methods.** We use squared exponential (SE) covariance functions for our experiments. The spectra of univariate SE covariance functions for all problems (see Appendix C.1) are determined using the Gaussian measure  $\mu = \mathcal{N}(0, 1)$ . The number of terms  $N_i$ ,  $i \in \{1, \dots, d\}$ , of each truncated univariate spectrum is determined such that  $\lambda_{i, N_i-1} / \lambda_{i, 1} \leq \eta_i$ , where  $\eta_i = 10^{-16}$ . If  $N_i > 1000$ , we set  $N_i = 1000$  to trade off between the accuracy of truncated spectra and computational cost. We also set  $n_o = 500$ . The maximum size of the exploration set is  $n_e = 250$ . The maximum size of the exploitation set is  $n_x = 200$ .

The number of initial observations is  $10d$  for all problems. The standard deviation of observation noise  $\sigma_n = 10^{-6}$  is applied for standardized output observations. The number of BO iterations for the 2D Schwefel and 4D Rosenbrock functions is 200, while that for the 10D Levy, 16D Ackley, and 16D Powell functions is 800. Other GP-TS methods for optimization of benchmark test functions including TS-DSRF (i.e., TS using decoupled sampling with random Fourier features) and TS-RF (i.e., TS using random Fourier features) are characterized by a total of 2000 random Fourier features.

To ensure a fair comparison of outer-optimization results, we first implement TS-roots and record the number of starting points used in each optimization iteration. We then apply other BO methods, each employing a gradient-based multistart optimizer with the same number of random starting points and identical termination criteria as those used for TS-roots in each iteration.

For the comparative inner-loop optimization performance of the proposed method via rootfinding with the random multistart and genetic algorithm approaches, we set the same termination tolerance on the objective function value as the stopping criterion for the methods. In addition, the number of starting points for the random multistart and the population size of the genetic algorithm are the same as the number of points in both the exploration and exploitation sets of rootfinding in each optimization iteration.

**Computational Tools.** We carry out all experiments, except those for inner-loop optimization, using a designated cluster at our host institution. This cluster hosts 9984 Intel CPU cores and 327680 Nvidia GPU cores integrated within 188 compute and 20 GPU nodes. The inner-loop optimization is implemented on a PC with an Intel® Core™ i7-1165G7 @ 2.80 GHz and 16 GB memory.

For the univariate global rootfinding via Chebyshev polynomials, we use MATLAB’s Chebfun package (Battles & Trefethen, 2004) and its corresponding implementation in Python, called chebpy (Richardson, 2016).

## I ADDITIONAL RESULTS

**Distance to Global Minimum.** Figure 9 shows the solution locations from 20 runs of TS-roots, TS-DSRF, TS-RF, EI, and LCB for the 2D Schwefel, 4D Rosenbrock, 10D Levy, 16D Ackley, 16D Powell functions.

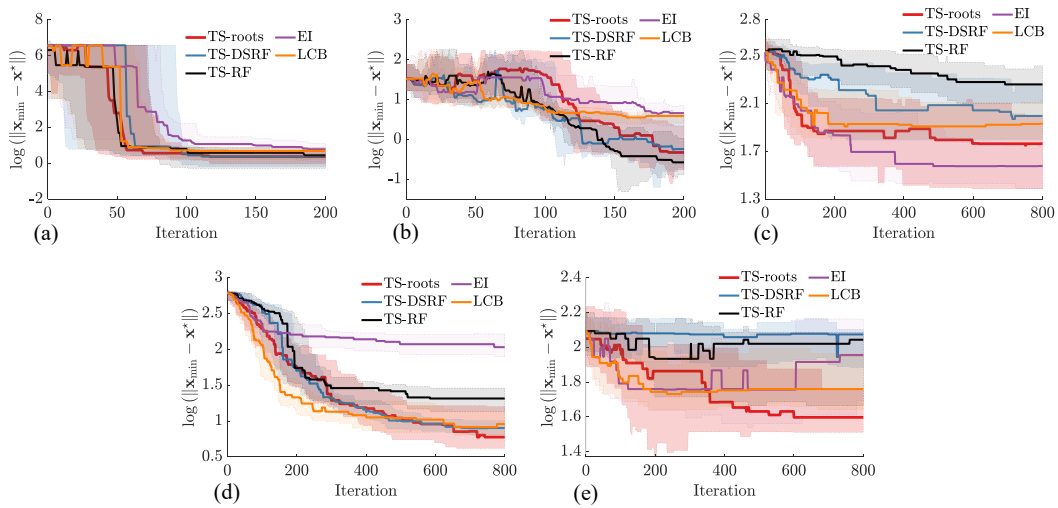


Figure 9: Outer-loop optimization results for (a) the 2D Schwefel, (b) 4D Rosenbrock, (c) 10D Levy, (d) 16D Ackley, (e) 16D Powell functions. The plots are histories of medians and interquartile ranges of solution locations from 20 runs of TS-roots, TS-DSRF (i.e., TS using decoupled sampling with random Fourier features), TS-RF (i.e., TS using random Fourier features), EI, and LCB.

**Comparison of Inner-loop Optimization Results.** Figures 10 and 11 compare the performance of the inner-loop optimization by three different initialization schemes, i.e., rootfinding, uniform grid, and Latin hypercube sampling, for low-dimensional cases of the 2D Schwefel and 4D Rosenbrock functions, and for higher-dimensional cases of the 10D Levy, 16D Ackley, and 16D Powell functions. Rootfinding performs better than the uniform grid and Latin hypercube sampling initialization schemes, especially in high-dimensional settings.

**Sample-average Posterior Function.** Figure 12 shows how we can improve the exploitation of GP-TS when increasing the exploration–exploitation control parameter  $N_c$ .

**Performance of Sample-average TS-roots.** Figure 13 shows the performance of sample-average TS-roots with different exploration–exploitation control parameters  $N_c$  for the 2D Schwefel, 4D Rosenbrock, and 6D Ackley functions.

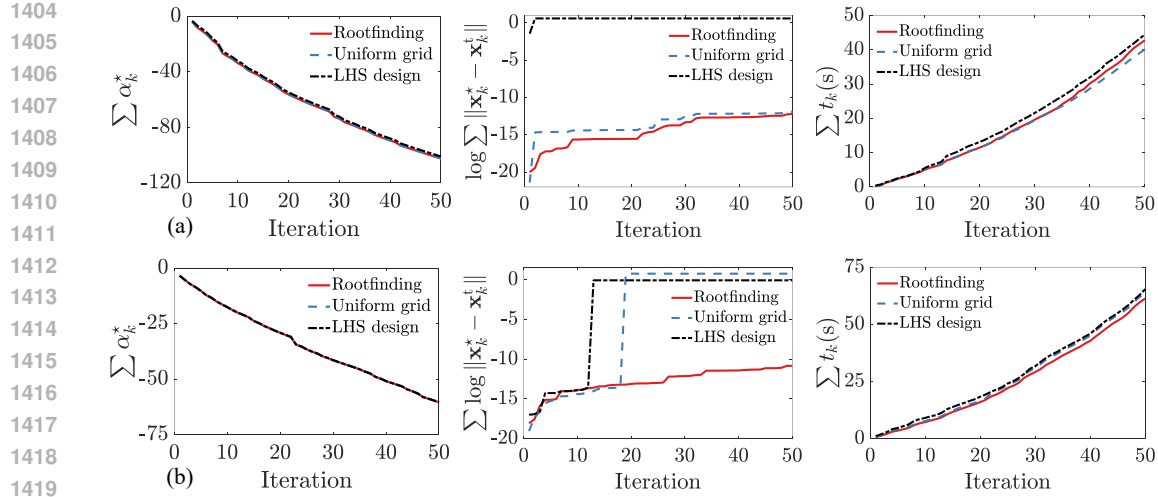


Figure 10: Inner-loop optimization results by three different initialization schemes, i.e., rootfinding, uniform grid, and Latin hypercube sampling, for (a) the 2D Schwefel and (b) 4D Rosenbrock functions. The plots are cumulative values of optimized GP-TS acquisition functions  $\alpha_k^*$ , cumulative distances between new solution points  $\mathbf{x}_k^*$  and the true global minima  $\mathbf{x}_k^t$  of the acquisition functions, and cumulative CPU times  $t_k$  for optimizing the acquisition functions.

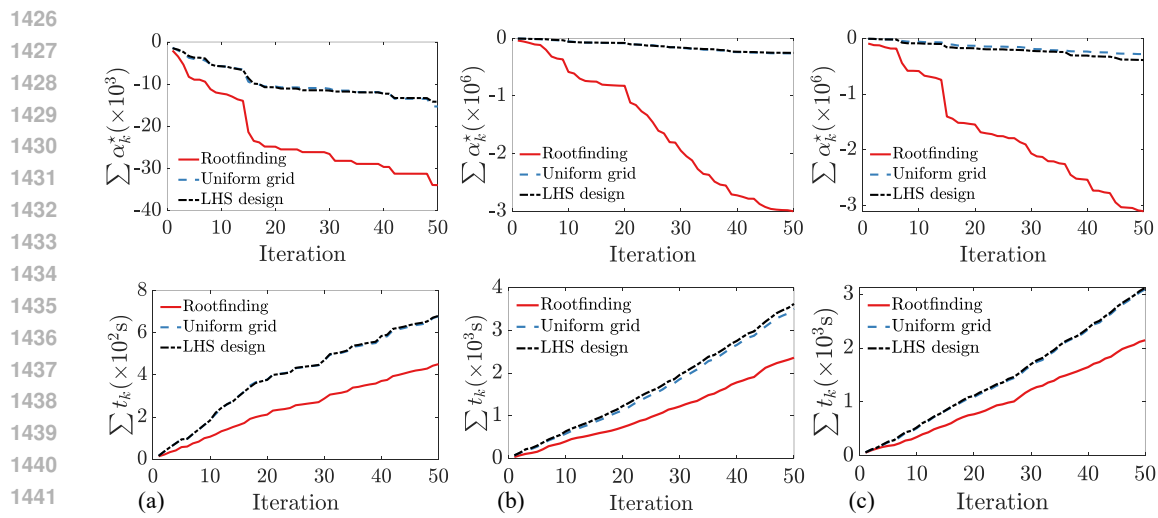


Figure 11: Inner-loop optimization results by three different initialization schemes, i.e., rootfinding, uniform grid, and Latin hypercube sampling, for (a) the 10D Levy, (b) 16D Ackley, and (c) 16D Powell functions. The plots are cumulative values of optimized GP-TS acquisition functions  $\alpha_k^*$  and cumulative CPU times  $t_k$  for optimizing the acquisition functions.

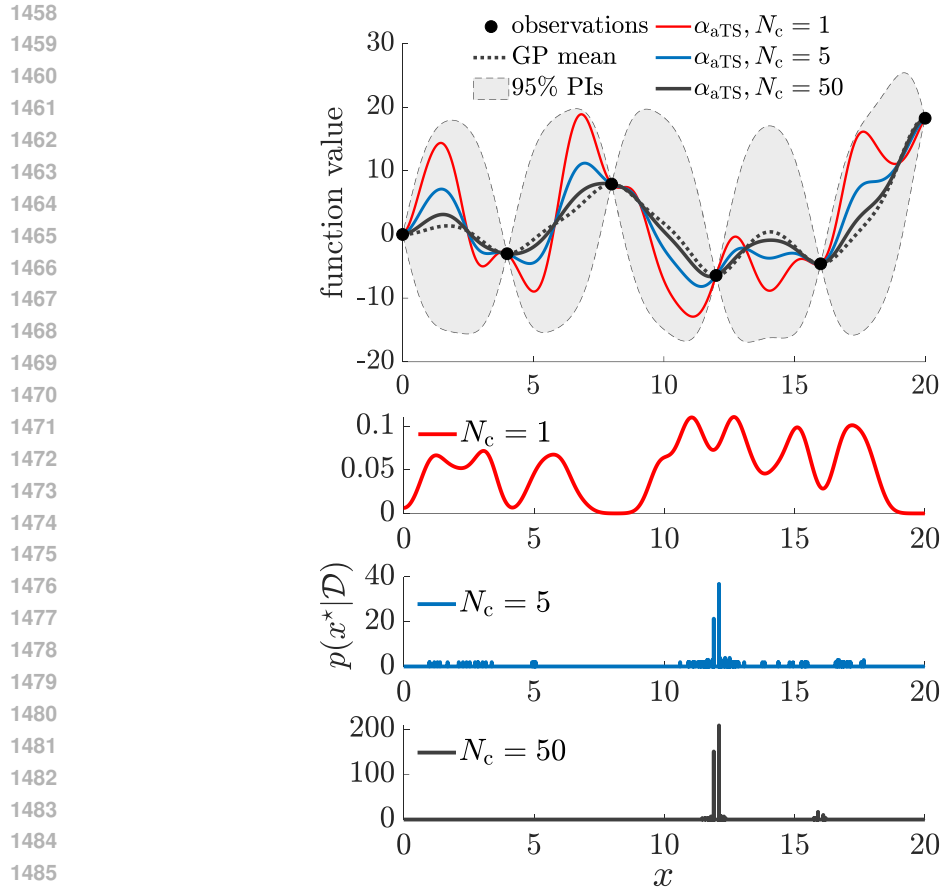


Figure 12: Sample-average posterior function for different values of  $N_c$ . The posterior function approaches the GP mean and the conditional distribution of the solution location  $p(x^*|\mathcal{D})$  is more concentrated when we increase  $N_c$ .

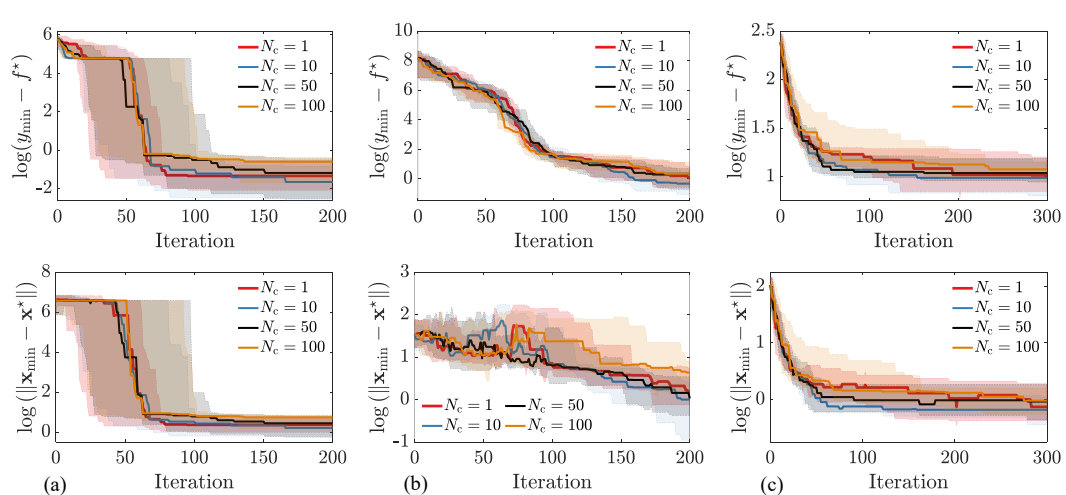


Figure 13: Performance of sample-average TS-roots with different control values  $N_c$  for (a) the 2D Schwefel, 4D Rosenbrock, and (b) 6D Ackley functions. The plots are histories of medians and interquartile ranges of solution values and solution locations from 20 runs of TS-roots for each  $N_c$  value.