

A dimensionality reduction technique for unconstrained global optimization of functions with low effective dimensionality

CORALIA CARTIS

Mathematical Institute, University of Oxford, Woodstock Road, Oxford OX2 6GG, UK, and The Alan Turing Institute, The British Library, London NW1 2DB, UK

Email: cartis@maths.ox.ac.uk

AND

ADILET OTEMISSOV[†]

The Alan Turing Institute, The British Library, London NW1 2DB, UK, and Mathematical Institute, University of Oxford, Woodstock Road, Oxford OX2 6GG, UK

[†]Corresponding author. Email: otemissoy@maths.ox.ac.uk, aotemissoy@nu.edu.kz

[Received on 16 January 2020; revised on 20 January 2021; accepted on 21 February 2021]

We investigate the unconstrained global optimization of functions with low effective dimensionality, which are constant along certain (unknown) linear subspaces. Extending the technique of random subspace embeddings in Wang *et al.* (2016, *J. Artificial Intelligence Res.*, **55**, 361–387), we study a generic Random Embeddings for Global Optimization (REGO) framework that is compatible with any global minimization algorithm. Instead of the original, potentially large-scale optimization problem, within REGO, a Gaussian random, low-dimensional problem with bound constraints is formulated and solved in a reduced space. We provide novel probabilistic bounds for the success of REGO in solving the original, low effective-dimensionality problem, which show its independence of the (potentially large) ambient dimension and its precise dependence on the dimensions of the effective and embedding subspaces. These results significantly improve existing theoretical analyses by providing the exact distribution of a reduced minimizer and its Euclidean norm and by the general assumptions required on the problem. We validate our theoretical findings by extensive numerical testing of REGO with three types of global optimization solvers, illustrating the improved scalability of REGO compared with the full-dimensional application of the respective solvers.

Keywords: global optimization; random matrix theory; dimensionality reduction techniques; functions with low effective dimensionality.

1. Introduction

In this paper, we address the unconstrained global optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^D} f(\mathbf{x}), \quad (\text{P})$$

where $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is a real-valued continuous, possibly non-convex, deterministic, function defined on the whole \mathbb{R}^D . We assume that there exists $\mathbf{x}^* \in \mathbb{R}^D$ such that $\min_{\mathbf{x} \in \mathbb{R}^D} f(\mathbf{x}) = f(\mathbf{x}^*) = f^*$. This implies that f is bounded below, namely $f^* > -\infty$, and that the minimum in (P) is attained (not all minimizers are at infinity).

To alleviate the curse of dimensionality, we further restrict ourselves to a particular class of functions whose true (intrinsic) dimension is much less than the ambient problem dimension. These functions are constant along certain linear subspaces, which may not necessarily be aligned with the standard axes. In literature, these functions are known under different names: functions with *low effective dimensionality* [49], functions with *active subspaces* [9] and *multi-ridge* functions [19, 46]. They have been found in a number of applications mainly related to parameter studies. In hyper-parameter optimization for neural networks [2] and heuristic algorithms for combinatorial optimization problems [26], studies have shown that the respective objective functions are affected by only a few hyper-parameters while the many other input hyper-parameters are redundant. Similarly, in complex engineering and physical simulation problems [9, 33], such as in climate modelling [29], systems are modelled by several input parameters with only a small number of the parameters or a combination of them having a true effect on the system's behaviour. To clarify this concept, we give a simple example of a function with lower effective dimensionality.

EXAMPLE 1.1 Consider the following optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) = \sin^2(x_1 - x_2 - 0.5). \quad (1.1)$$

By solving $f(\mathbf{x}) = 0$, we find that the set of global minimizers is given by $\{(1 \ -1)^T t - (0 \ 0.5 + \pi k)^T : t \in \mathbb{R}, k \in \mathbb{Z}\}$. For each fixed value of k , the set corresponds to a distinct line of global minimizers along which the function is constant (Fig. 1). The effective subspace¹ of f is $(x_1 \ x_2) = (1 \ -1)^T y$ for $y \in \mathbb{R}$. We substitute this in (1.1) to obtain the reduced/lower-dimensional optimization problem $\min_{y \in \mathbb{R}} \sin^2(2y - 0.5)$, which has the same global minimum as (1.1), with global minimizers $y_k^* = \pi k/2 + 0.25$, $k \in \mathbb{Z}$. We recover the corresponding solutions to (1.1) by setting $\mathbf{x}_k^* = (1 \ -1)^T y_k^*$ for $k \in \mathbb{Z}$.

As Example 1.1 illustrates, it is possible to cast (P) into a lower-dimensional problem which has the same global minimum f^* . This is straightforward when the effective subspace is known but far less so in applications where f is potentially black-box. When the effective subspace is unknown, it was proposed (in the context of Bayesian optimization (BO)) in Wang *et al.* [49] to use random embeddings. The proposed technique solves the following lower-dimensional optimization problem instead of directly tackling (P):

$$\begin{aligned} \min_{\mathbf{y}} f(\mathbf{A}\mathbf{y}) \\ \text{subject to } \mathbf{y} \in \mathcal{Y} = [-\delta, \delta]^d, \end{aligned} \quad (\text{RP})$$

where \mathbf{A} is a $D \times d$ Gaussian random matrix (see Definition A.1) and $\mathcal{Y} = [-\delta, \delta]^d$ for some carefully chosen $\delta > 0$ and $d \ll D$. Note that, unlike (P), (RP) has (box) constraints, which are typically imposed to make the approach practical (i.e. to avoid unrealistic searches over infinite domains).

DEFINITION 1.2 We say that (RP) is *successful* if there exists $\mathbf{y}^* \in \mathcal{Y}$ such that $f(\mathbf{A}\mathbf{y}^*) = f^*$.

¹ The effective subspace can be determined by considering the orthogonal complement of the constant subspace (along which f does not vary), in this example, spanned by the vector $(1 \ 1)^T$.

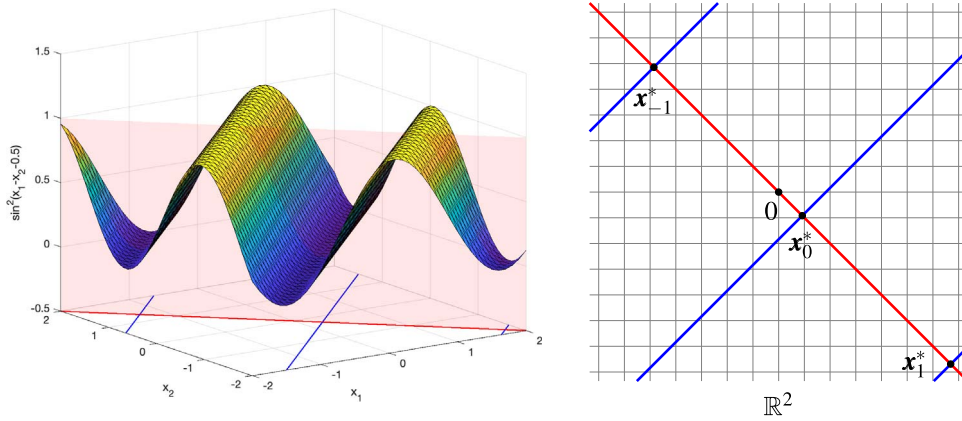


FIG. 1. The function in (1.1) and its domain are plotted on the left and right sides, respectively. The red line is the effective subspace $x = (1 - 1)y$ and it intersects the blue lines of global minimizers at x_k^* defined in Example 1.1 for $k = -1, 0, 1$; these points also correspond to optimal solutions in the reduced space.

Related work. The scalability challenges of BO algorithms for generic black-box functions have prompted research into improving efficiency of this class of methods for functions with special structure. Different structural assumptions on the objective have been analysed for BO, such as additivity or (partial) separability, which assumes that the objective function can be represented as the sum of smaller-dimensional functions with non-overlapping variables [28, 32, 48] or with overlapping ones [40].

Another popular structural assumption is the above-mentioned low-effective dimensionality of the objective. In its simplest form, this considers the effective subspace to be aligned with the coordinate axes, which is equivalent to the presence of inactive variables [1, 8]. More generally, the optimization of functions that are constant along *arbitrary* linear subspaces—which, as mentioned above, is also the focus of this paper—has been addressed using BO methods in [13, 15, 22, 49] and extended to other problem and algorithm classes such as derivative-free optimization [38], evolutionary [43] and genetic [11] methods and multi-objective optimization [39]. Some proposals learn the effective subspace of the function beforehand (using, for example, a low-rank matrix recovery approach) [19, 46] and then optimize in the reduced subspace [13, 15]. Alternating learning and optimization steps has also been proposed [22], as well as bypassing learning and directly optimizing in randomly chosen low-dimensional subspaces (provided an estimate of the effective dimension is known) [4, 5, 49].

For the latter, Wang *et al.* [49] developed the so-called REMBO algorithm, which is a BO framework for problem (P) with box constraints $x \in \mathcal{X}$ that uses Gaussian random embeddings (namely, A is a Gaussian random matrix) to generate the reduced problem (RP). They find that the size of \mathcal{Y} is the primary factor in determining the success (or failure) of the reduced problem and quantify the probability of success of (RP) for the case when the embedded dimension d is equal to the effective one and the effective subspace is aligned with the coordinate axes [49, Theorem 3]. A challenge of (RP) for BO with box constraints is that, even when (RP) is successful, the high-dimensional image $Ay \in \mathbb{R}^D$ of a point $y \in \mathcal{Y}$ may be outside the feasible set \mathcal{X} . For this reason, REMBO is equipped with a map $p_{\mathcal{X}} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ that projects the image of the reduced solutions that fall outside \mathcal{X} to the closest point on the boundary of \mathcal{X} . To model a Gaussian process for the reduced problem, [49] proposes two kernels: a high-dimensional $k_{\mathcal{X}}$ and a low-dimensional $k_{\mathcal{Y}}$. Kernel $k_{\mathcal{X}}$ suffers from high-dimensionality as it constructs a GP in a D -dimensional space. The benefit of $k_{\mathcal{Y}}$ is that it constructs a GP in a d -dimensional

subspace, but this kernel over-explores regions in \mathcal{Y} whose high-dimensional images outside \mathcal{X} are mapped to the same points in \mathcal{X} through a non-injective $p_{\mathcal{X}}$. To remedy these issues, Binois *et al.* [4] propose a new kernel $k_{\mathcal{Y}}$ which has the benefit of being low-dimensional while avoiding the over-exploratory tendency of $k_{\mathcal{X}}$. Binois *et al.* [5] also propose a new mapping γ (instead of $p_{\mathcal{X}}$) and define \mathcal{Y} and new kernels based on this new mapping.

Sangyang & Kabán [43] develop REMEDA, which uses random embeddings within an evolutionary algorithm EDA. Their theoretical results on quantifying the size of \mathcal{Y} / the success of (RP) improve on those in [49] and are applicable for certain choices of d , greater than the effective subspace dimension; they also experiment with estimating the effective dimension numerically.

Qian *et al.* [38] extend the framework and some of the results in Wang *et al.* [49] to functions with approximate low effective subspaces, proposing the use of multiple random embeddings within any derivative-free solver. They contrast the use of a single versus multiple embeddings on three test problems of varying dimensions and using three different types of derivative-free solvers (evolutionary, Bayesian and model based).

Recently, in the context of BO, Nayebi *et al.* [34] use a different random ensemble based on hashing matrices to represent the embedded subspaces and define \mathcal{Y} as $[-1, 1]^d$; this formulation guarantees that the high-dimensional points are always inside \mathcal{X} and, thus, their method avoids the feasibility corrections of REMBO.

Our contributions. We investigate a general random embeddings framework for unconstrained global optimization of functions with low effective dimensionality, where we allow the effective subspace of the objective function and its dimension (denoted by d_e) to be arbitrary (not necessarily aligned with coordinate axes and not limited in dimension by problem constants). This framework also allows the use of any global solver to solve the reduced problem. We note that, as problem (P) has no (bound) constraints, we do not need to use the projection operator $p_{\mathcal{X}}$ in [4, 5, 49]. Of course, this comes at the cost of our approach being unable to guarantee feasibility for the original problem if (P) does have constraints.

We significantly extend and improve the theoretical analyses in [43, 49], providing an in-depth investigation of the reduced problem (RP) when Gaussian random embeddings are used. In particular, while [43, 49] estimate the Euclidean norm of a (random) reduced minimizer, we derive its exact distribution, using tools from random matrix theory. We show that this reduced minimizer, when appropriately scaled, follows the inverse chi-squared distribution with $d - d_e + 1$ degrees of freedom, where d is the dimension of the random embedding (Theorem 3.7). Moreover, we derive the probability density function of this reduced minimizer (Theorem 3.10) by first proving that it follows a spherical distribution. These results imply that, under certain assumptions, solving (RP) has no dependence on the ambient dimension D . Subsequently, Theorem 4.1 and Corollary 4.2 estimate the probability that (RP) is successful. The latter result extends both [49, Theorem 3] and [43, Theorem 2] to arbitrary effective subspaces and any $d \geq d_e$ and establishes a notable and more precise trade-off between the success of (RP), δ (the size of the reduced domain \mathcal{Y}) and the embedding dimension d , thus allowing us to choose appropriate values for these parameters in the algorithm. Furthermore, we describe how to extend the main results to affine random embeddings (which draw random subspaces at any chosen (reference) point in \mathbb{R}^D), which indicate that the probability of success of (RP) is higher if the point of reference is closer to the set of global minimizers.

Similarly to the algorithmic frameworks proposed in [38, 49], we propose Random Embeddings for Global Optimization (REGO) that solves a single randomly embedded reduced problem (RP) instead of

(P) and is compatible with any generic global optimization solver². We use and validate our theoretical results by providing extensive numerical testing of REGO with three types of solvers for (RP): DIRECT (Lipschitz-optimization), BARON (branch and bound) and KNITRO (multi-start local optimization). We use 19 standard global optimization test problems to generate functions with effective dimensionality structure and of growing ambient dimension D . When comparing REGO with the direct optimization of the ensuing problems without embeddings, we find that REGO's performance is essentially independent of D for all three solvers and that it is successful in recovering the original global minimum in most cases with only one embedding³. We also test the robustness of REGO's performance to variations in algorithm parameters such as δ and d .

Paper outline. In Section 2, we formally define and describe functions with low effective dimensionality emphasizing their geometrical aspects. In Section 3, we characterize the reduced minimizers in the reduced space focusing on the minimal 2-norm minimizer. For this minimizer, we derive the distribution of its Euclidean norm and its probability density function. We use the former result in Section 4 to derive a probabilistic bound for the success of (RP). In Section 5, we conduct numerical experiments to test REGO algorithm on functions with low effective dimensionality using three optimization solvers, namely DIRECT, BARON and KNITRO, while in Section 6, we draw our conclusions and future directions.

Notation. We use bold capital letters to denote matrices (\mathbf{A}) and bold lowercase letters (\mathbf{a}) to denote vectors. In particular, we use \mathbf{I}_D to denote the $D \times D$ identity matrix and $\mathbf{0}_D$, $\mathbf{1}_D$ (or simply $\mathbf{0}$, $\mathbf{1}$) to denote the D -dimensional vectors of zeros and ones, respectively. For an $D \times d$ matrix \mathbf{A} , we write $\text{range}(\mathbf{A})$ to denote the linear subspace spanned by the columns of \mathbf{A} in \mathbb{R}^D .

We let $\langle \cdot, \cdot \rangle$, $\| \cdot \|$ and $\| \cdot \|_\infty$ denote the usual Euclidean inner product, the Euclidean norm and the infinity norm, respectively. Where emphasis is needed, for the Euclidean norm, we also use $\| \cdot \|_2$.

Given two random variables (vectors) x and y (\mathbf{x} and \mathbf{y}), we write $x \stackrel{\text{law}}{=} y$ ($\mathbf{x} \stackrel{\text{law}}{=} \mathbf{y}$) to denote the fact that x and y (\mathbf{x} and \mathbf{y}) have the same distribution. We reserve the letter \mathbf{A} to refer to a $D \times d$ Gaussian random matrix (see Definition A.1) and write χ_n^2 to denote a chi-squared random variable with n degrees of freedom (see Definition A.2).

2. Functions with low effective dimensionality

In this section, we formally define functions with low effective dimensionality and describe the geometry of (RP).

2.1 Definitions and assumptions

Functions with low effective dimensionality can be defined in at least two ways [19, 49]. We will work with a definition given in terms of linear subspaces, provided in [49].

DEFINITION 2.1 (Functions with low effective dimensionality). A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ has effective dimensionality $d_e \leq D$ if there exists a linear subspace \mathcal{T} of dimension d_e such that for all vectors \mathbf{x}_\top

² The REGO framework is defined for the unconstrained (P) but can also be helpful for constrained problems (for example, $\mathbf{x} \in \mathcal{X}$), where the constraints are imposed just to avoid searches over an infinite domain and where minimizers outside the feasible domain are acceptable.

³ These numerical results assume that (an upper bound on) the true effective dimension d_e is known/available.

in \mathcal{T} and \mathbf{x}_\perp in \mathcal{T}^\perp (orthogonal complement of \mathcal{T}), we have

$$f(\mathbf{x}_\top + \mathbf{x}_\perp) = f(\mathbf{x}_\top), \quad (2.1)$$

and d_e is the smallest integer satisfying (2.1).

The linear subspace \mathcal{T} is called the *effective* subspace of f and its orthogonal complement \mathcal{T}^\perp , the *constant* subspace of f . It is convenient to think of \mathcal{T}^\perp as a subspace of no variation of largest dimension (along which the value of f does not change) and \mathcal{T} as its orthogonal complement.

Every vector \mathbf{x} can be decomposed as $\mathbf{x} = \mathbf{x}_\top + \mathbf{x}_\perp$, where \mathbf{x}_\top and \mathbf{x}_\perp are orthogonal projections of \mathbf{x} onto \mathcal{T} and \mathcal{T}^\perp , respectively. In particular, if \mathbf{x}^* is a global minimizer and f^* is the global minimum of f in \mathcal{X} then $\mathbf{x}^* = \mathbf{x}_\top^* + \mathbf{x}_\perp^*$ and

$$f^* = f(\mathbf{x}^*) = f(\mathbf{x}_\top^* + \mathbf{x}_\perp^*) = f(\mathbf{x}_\top^*). \quad (2.2)$$

Moreover, we have

$$f^* = f(\mathbf{x}_\top^*) = f(\mathbf{x}_\top^* + \mathbf{x}_\perp)$$

for every vector \mathbf{x}_\perp in \mathcal{T}^\perp . It is important to note that there can be multiple points \mathbf{x}_\top^* in \mathbb{R}^D satisfying the above definition such as, for instance, \mathbf{x}_{-1}^* , \mathbf{x}_0^* and \mathbf{x}_1^* in Example 1.1. By contrast, the function $f = (x_1 - x_2 - 0.5)^2$ admits a unique \mathbf{x}_\top^* given by $(0.25 \ -0.25)^T$.

We summarize the above discussion in the following assumption.

ASSUMPTION 2.2 The function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is continuous and has effective dimensionality $d_e \leq d$ with effective subspace⁴ \mathcal{T} and constant subspace \mathcal{T}^\perp spanned by the columns of orthonormal matrices $\mathbf{U} \in \mathbb{R}^{D \times d_e}$ and $\mathbf{V} \in \mathbb{R}^{D \times (D-d_e)}$, respectively.

Recalling the definition of problem (P) on page 6, let

$$\mathcal{G} = \{\mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) = f^*\}$$

be the set of global minimizers in \mathbb{R}^D . Under Assumption 2.2, the set \mathcal{G} can be represented as a union of (possibly infinitely many) affine subspaces each containing one particular \mathbf{x}_\top^* (see proof of Theorem 4.1). Each of these affine subspaces is a $(D - d_e)$ -dimensional set $\{\mathbf{x} \in \mathbb{R}^D : \mathbf{x} \in \mathbf{x}_\top^* + \mathcal{T}^\perp\}$ —a translation of \mathcal{T}^\perp by the vector \mathbf{x}_\top^* that the corresponding affine subspace must contain. In particular, if there is a unique \mathbf{x}_\top^* in \mathbb{R}^D , then $\mathcal{G} = \{\mathbf{x} \in \mathbb{R}^D : \mathbf{x} \in \mathbf{x}_\top^* + \mathcal{T}^\perp\}$. Note also that point(s) \mathbf{x}_\top^* lie in $\mathcal{G} \cap \mathcal{T}$ and are the closest minimizers to the origin in Euclidean norm among all the minimizers lying in their respective affine subspaces.

Our analysis applies to any minimizer \mathbf{x}^* with $\mathbf{x}_\top^* \neq \mathbf{0}$. If $\mathbf{x}_\top^* = \mathbf{0}$, then (RP) has a trivial solution. In that case, the origin is a global minimizer implying that every embedding is successful with a solution $\mathbf{y}^* = \mathbf{0}$. Hence, we focus our analysis for finding a(ny) minimizer $\mathbf{x}^* \in \mathcal{G}$ with $\mathbf{x}_\top^* \neq \mathbf{0}$.

⁴ Note that \mathcal{T} in Assumption 2.2 may not be aligned with the standard axes.

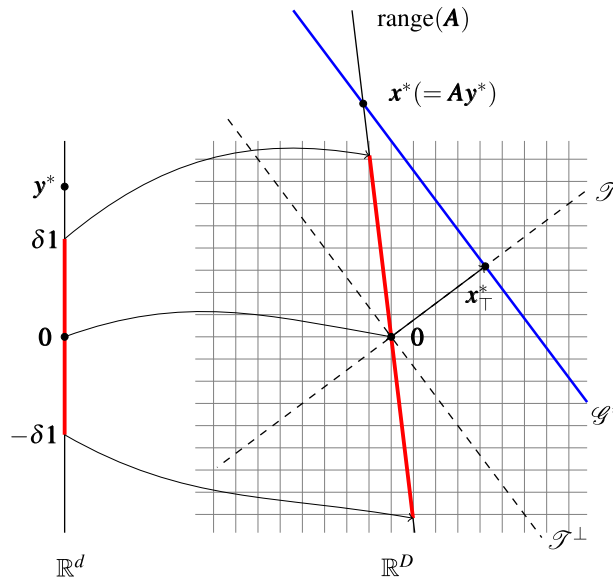


Fig. 2. The figure shows an abstract illustration of the embedding of a d -dimensional linear subspace into \mathbb{R}^D . The line $\text{range}(\mathbf{A})$ corresponds to the embedded subspace. The red line in \mathbb{R}^d represents the hypercube $\mathcal{Y} = \{\mathbf{y} \in \mathbb{R}^d : -\delta \mathbf{1} \leq \mathbf{y} \leq \delta \mathbf{1}\}$, which, after application of \mathbf{A} , maps to the red line along $\text{range}(\mathbf{A})$ in \mathbb{R}^D . In this configuration, condition (2.4) is satisfied but (2.3) is not: $\text{range}(\mathbf{A})$ intersects \mathcal{G} at $\mathbf{x}^* = \mathbf{A}\mathbf{y}^*$ but \mathbf{y}^* lies outside \mathcal{Y} .

ASSUMPTION 2.3 Given Assumption 2.2, let $\mathbf{x}^* \in \mathcal{G}$ such that $\mathbf{x}_\top^* = \mathbf{U}\mathbf{U}^T\mathbf{x}^*$ —the unique Euclidean projection of \mathbf{x}^* onto \mathcal{T} —is non-zero. Let $\mathcal{G}^* := \mathbf{x}_\top^* + \mathcal{T}^\perp$ be the affine subspace of \mathcal{G} that contains \mathbf{x}_\top^* .

The set \mathcal{G} contains infinitely many global minimizers—a particularly useful feature of the functions with low effective dimensionality; this fact suggests that targeting \mathcal{G} numerically may potentially be easier than if \mathcal{G} contained only one point.

2.2 Geometric description

We now provide a geometric description of (RP), which serves as a basis for our theoretical investigations.

In Fig. 2, we illustrate schematically \mathcal{T} (the effective subspace of f), \mathcal{T}^\perp (the orthogonal component of \mathcal{T}), \mathcal{G}^* (a connected component⁵ of \mathcal{G}) and \mathbf{x}_\top^* (the orthogonal projection of the global minimizers on \mathcal{G}^* onto \mathcal{T}). Since the orientation and position of these geometric objects are solely determined by the (deterministic) objective function, they are fixed, non-random.

By applying the ‘random embedding’ (RP), we switch from optimizing over \mathbb{R}^D to optimizing over \mathcal{Y} . The linear mapping $\mathbf{y} \rightarrow \mathbf{A}\mathbf{y}$ maps points of the hypercube \mathcal{Y} to points along the subspace $\text{range}(\mathbf{A})$ in \mathbb{R}^D , which means that searching over \mathcal{Y} is equivalent to searching over the corresponding feasible set along $\text{range}(\mathbf{A})$ in \mathbb{R}^D . An example of this mapping is illustrated in Fig. 2 with two red line segments: the segment (from $-\delta\mathbf{1}$ to $\delta\mathbf{1}$) representing \mathcal{Y} is being mapped to the right segment, which lies in $\text{range}(\mathbf{A})$.

⁵ Recall that \mathcal{G} is a union of affine subspaces; \mathcal{G}^* is one of them.

It is important to note that the centre of \mathcal{Y} maps to the origin in \mathbb{R}^D and, hence, the corresponding search in the original space is also centred at the origin.

The most valuable information that we want to retain while performing dimensionality reduction is the value of f^* . We would like $\min_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{A}\mathbf{y}) = f^*$, which holds only if there is at least one \mathbf{y}^* in \mathcal{Y} such that $f(\mathbf{A}\mathbf{y}^*) = f^*$. This condition has a geometric interpretation and, following from the definition of \mathcal{G} , it can equivalently be stated as

$$\text{there exists a } \mathbf{y}^* \in \mathcal{Y} \text{ such that } \mathbf{A}\mathbf{y}^* \in \mathcal{G}. \quad (2.3)$$

For (2.3) to hold, we must first ensure that

$$\text{there exists a } \mathbf{y}^* \text{ in } \mathbb{R}^d \text{ such that } \mathbf{A}\mathbf{y}^* \in \mathcal{G}. \quad (2.4)$$

In this regard, Wang *et al.* [49] proved the following theorem.

THEOREM 2.4 (see [49, Theorem 2]) Let Assumption 2.2 hold, and let \mathbf{A} be a $D \times d$ Gaussian matrix with $d \geq d_e$. Then, with probability one, for any fixed $\mathbf{x} \in \mathbb{R}^D$, there exists a $\mathbf{y} \in \mathbb{R}^d$ such that $f(\mathbf{x}) = f(\mathbf{A}\mathbf{y})$. In particular, for a global minimizer \mathbf{x}^* , with probability one, there exists a $\mathbf{y}^* \in \mathbb{R}^d$ such that $f(\mathbf{A}\mathbf{y}^*) = f(\mathbf{x}^*) = f^*$.

While satisfaction of (2.4) only depends on $d \geq d_e$, that of (2.3) is determined by the values of both d and δ . For larger values of d and/or δ the probability that (2.3) is satisfied is higher. One must, on the other hand, be cognisant of the fact that larger values of d —the dimension of (RP)—and/or δ (the half-length of the domain) demand more computational resources. Therefore, a careful calibration of these two parameters is needed to ensure that (RP) is successful for most embeddings, at the same time being capable to converge to the solution within the computational budget. In this regard, our analysis will attempt to answer the following question: what are optimal values of d and δ such that (2.3) is satisfied with ‘high’ probability?

3. Characterizing minimizers in the reduced space

The analysis of this section focuses on determining the distribution of the random minimizer \mathbf{y}^* of $f(\mathbf{A}\mathbf{y})$, which satisfies $f(\mathbf{A}\mathbf{y}^*) = f^*$. These results will inform us on the effects of the different parameters on the success of (RP) allowing us to estimate the values of δ and d that are likely to increase the chances of successful recovery of f^* .

The following theorem provides a useful characterization of \mathbf{y}^* . The theorem and its proof were inspired by the proofs of [49, Theorems 2 and 3].

THEOREM 3.1 Let Assumption 2.2 hold, and let \mathbf{x}_\top^* and \mathcal{G}^* be defined as in Assumption 2.3. Let \mathbf{A} be a $D \times d$ Gaussian matrix. Then, $\mathbf{y}^* \in \mathbb{R}^d$ satisfies $\mathbf{A}\mathbf{y}^* \in \mathcal{G}^*$ if and only if

$$\mathbf{B}\mathbf{y}^* = \mathbf{z}^*, \quad (3.1)$$

where the $d_e \times d$ random matrix \mathbf{B} satisfies $\mathbf{B} = \mathbf{U}^T \mathbf{A}$ and where $\mathbf{z}^* \in \mathbb{R}^{d_e}$ is uniquely defined by $\mathbf{U}\mathbf{z}^* = \mathbf{x}_\top^*$.

Proof. Let $\mathbf{y}^* \in \mathbb{R}^d$ be such that $\mathbf{A}\mathbf{y}^* \in \mathcal{G}^*$. First, we establish that

$$\mathbf{A}\mathbf{y}^* \in \mathcal{G}^* \text{ if and only if } \mathbf{x}_\top^* = \mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^*. \quad (3.2)$$

Suppose that $\mathbf{A}\mathbf{y}^* \in \mathcal{G}^*$. Then, using the definition of \mathcal{G}^* in Assumption 2.3, we can write $\mathbf{A}\mathbf{y}^* = \mathbf{x}_\top^* + \mathbf{x}_\perp$ for some $\mathbf{x}_\perp \in \mathcal{T}^\perp$. The orthogonal projection of $\mathbf{A}\mathbf{y}^*$ onto \mathcal{T} is given by

$$\mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^* = \mathbf{U}\mathbf{U}^T(\mathbf{x}_\top^* + \mathbf{x}_\perp) = \mathbf{x}_\top^*,$$

where we have used $\mathbf{U}\mathbf{U}^T\mathbf{x}_\top^* = \mathbf{x}_\top^*$ and $\mathbf{U}\mathbf{U}^T\mathbf{x}_\perp = \mathbf{0}$.

Conversely, assume that \mathbf{y}^* satisfies

$$\mathbf{x}_\top^* = \mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^*. \quad (3.3)$$

Denote by \mathbf{S} the $D \times D$ orthogonal matrix $(\mathbf{U} \ \mathbf{V})$, where \mathbf{V} is defined in Assumption 2.2. Using (3.3) and the identity $\mathbf{U}\mathbf{U}^T + \mathbf{V}\mathbf{V}^T = \mathbf{S}\mathbf{S}^T = \mathbf{I}_D$, we obtain

$$\mathbf{A}\mathbf{y}^* = (\mathbf{U}\mathbf{U}^T + \mathbf{V}\mathbf{V}^T)\mathbf{A}\mathbf{y}^* = \mathbf{x}_\top^* + \mathbf{V}\mathbf{V}^T\mathbf{A}\mathbf{y}^*.$$

Note that $\mathbf{V}\mathbf{V}^T\mathbf{A}\mathbf{y}^*$ lies on \mathcal{T}^\perp as it is the orthogonal projection of $\mathbf{A}\mathbf{y}^*$ onto \mathcal{T}^\perp , which implies that $\mathbf{A}\mathbf{y}^* \in \mathcal{G}^*$. This completes the proof of (3.2).

Now, we show that (3.1) and (3.3) are equivalent. We multiply both sides of $\mathbf{x}_\top^* = \mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^*$ by \mathbf{S}^T and obtain

$$\begin{pmatrix} \mathbf{U}^T \\ \mathbf{V}^T \end{pmatrix} \mathbf{x}_\top^* = \begin{pmatrix} \mathbf{U}^T \\ \mathbf{V}^T \end{pmatrix} \mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^*. \quad (3.4)$$

Since \mathbf{x}_\top^* is in the column span of \mathbf{U} , it can be written as $\mathbf{x}_\top^* = \mathbf{U}\mathbf{z}^*$ for some (unique) vector $\mathbf{z}^* \in \mathbb{R}^{d_e}$. By substituting the above into (3.4), we obtain

$$\begin{pmatrix} \mathbf{U}^T\mathbf{U}\mathbf{z}^* \\ \mathbf{V}^T\mathbf{U}\mathbf{z}^* \end{pmatrix} = \begin{pmatrix} \mathbf{U}^T\mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^* \\ \mathbf{V}^T\mathbf{U}\mathbf{U}^T\mathbf{A}\mathbf{y}^* \end{pmatrix}.$$

This reduces to

$$\begin{pmatrix} \mathbf{z}^* \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{U}^T\mathbf{A}\mathbf{y}^* \\ \mathbf{0} \end{pmatrix},$$

where we have used the identities $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ and $\mathbf{V}^T\mathbf{U} = \mathbf{0}$, which follow from Assumption 2.2. To obtain (3.3) from (3.1), multiply (3.1) by \mathbf{U} . \square

REMARK 3.2 Thereafter, we write \mathbf{B} to refer to the $d_e \times d$ Gaussian matrix⁶ $\mathbf{U}^T\mathbf{A}$, where \mathbf{A} is a $D \times d$ Gaussian matrix and where \mathbf{U} is defined in Assumption 2.2. Furthermore, we write \mathbf{z}^* to refer to the

⁶ Since \mathbf{U} is orthogonal, it follows from Theorem A.2 that $\mathbf{B} = \mathbf{U}^T\mathbf{A}$ is a Gaussian matrix.

$d_e \times 1$ vector that satisfies $U\mathbf{z}^* = \mathbf{x}_\top^*$, where \mathbf{x}_\top^* is defined in Assumption 2.3. Observe that $\|\mathbf{z}^*\| = \|\mathbf{x}_\top^*\|$ since $U\mathbf{z}^* = \mathbf{x}_\top^*$ and U is orthogonal.

COROLLARY 3.3 Let Assumption 2.2 hold. Let $S^* = \{\mathbf{y}^* : \mathbf{A}\mathbf{y}^* \in \mathcal{G}^*\}$, where \mathbf{A} is a $D \times d$ Gaussian matrix and where \mathcal{G}^* is defined as in Assumption 2.3. Then, the following holds:

- if $d = d_e$, then S^* has exactly one element with probability 1;
- if $d > d_e$, then S^* has infinitely many elements with probability 1.

Proof. It follows from Theorem 3.1 that the set S^* and the set of solutions to $\mathbf{B}\mathbf{y} = \mathbf{z}^*$ coincide. According to Theorem A.3, $\mathbf{B}\mathbf{B}^T$ is positive definite with probability 1, which implies that $\text{rank}(\mathbf{B}\mathbf{B}^T) = d_e$ with probability 1. Since $\text{rank}(\mathbf{B}) = \text{rank}(\mathbf{B}\mathbf{B}^T)$, $\text{rank}(\mathbf{B}) = d_e$ with probability 1. Hence, the linear system $\mathbf{B}\mathbf{y} = \mathbf{z}^*$ almost surely has a solution. If $d = d_e$, the linear system has only one solution. If $d > d_e$ the system is underdetermined and, therefore, has infinitely many solutions. \square

3.1 Choosing a suitable minimizer

While S^* contains infinitely many solutions if $d > d_e$, it is sufficient that one of these solutions is contained in \mathcal{Y} for (RP) to be successful. We proceed further by choosing one particular solution \mathbf{y}^* that is easy to analyse and based on the analysis will adjust parameters δ and d appropriately to ensure that the chosen \mathbf{y}^* falls inside the feasible set \mathcal{Y} with high probability. The solutions that are likely to fall inside the feasible domain must be close to the origin. In this regard, we propose two candidates:

$$\begin{aligned} \mathbf{y}_2^* &= \underset{\mathbf{y} \in \mathbb{R}^d}{\text{argmin}} \|\mathbf{y}\|_2 \\ \text{s.t. } \mathbf{y} &\in S^*, \end{aligned} \quad (3.5)$$

$$\begin{aligned} \mathbf{y}_\infty^* &= \underset{\mathbf{y} \in \mathbb{R}^d}{\text{argmin}} \|\mathbf{y}\|_\infty \\ \text{s.t. } \mathbf{y} &\in S^*. \end{aligned} \quad (3.6)$$

Due to the definition of \mathcal{Y} as a box, the minimizer (3.6) with the minimal infinity norm is particularly of interest. Since \mathbf{y}_∞^* has the smallest infinity norm among all solutions in S^* , knowledge of \mathbf{y}_∞^* would allow us to choose the smallest possible \mathcal{Y} while ensuring that (RP) is successful. Despite this convenient fact, we found that it is more difficult to study \mathbf{y}_∞^* and have decided to investigate \mathbf{y}_2^* instead.

REMARK 3.4 For $d = d_e$, $\mathbf{y}_2^* = \mathbf{y}_\infty^*$ because S^* contains only one element.

LEMMA 3.5 Let Assumption 2.2 hold, and let \mathcal{G}^* and \mathbf{y}_2^* be defined as in Assumption 2.3 and (3.5), respectively. Problem (RP) is successful in the sense of Definition 1.2 if $\mathbf{y}_2^* \in \mathcal{Y}$.

Proof. Assume that $\mathbf{y}_2^* \in \mathcal{Y}$. Then, \mathbf{y}_2^* is a feasible solution of (RP). By the definitions of \mathbf{y}_2^* and S^* , we also have $\mathbf{A}\mathbf{y}_2^* \in \mathcal{G}^*$; this implies that $f(\mathbf{A}\mathbf{y}_2^*) = f^*$. Hence, (RP) is successful by Definition 1.2. \square

COROLLARY 3.6 Let Assumption 2.2 hold. Problem (3.5) has a unique solution given by

$$\mathbf{y}_2^* = \mathbf{B}^T (\mathbf{B}\mathbf{B}^T)^{-1} \mathbf{z}^*. \quad (3.7)$$

Proof. It follows from Theorem 3.1 that the solution(s) of (3.5) must be equal to the solution(s) of the following problem:

$$\begin{aligned} \min \quad & \|y\|_2^2 \\ \text{s.t.} \quad & By = z^*, \end{aligned}$$

which has the solution (3.7). \square

3.2 Distribution of minimal Euclidean norm minimizer

The present section derives the distribution of y_2^* and its Euclidean norm.

3.2.1 The distribution of the Euclidean norm of y_2^*

THEOREM 3.7 Let Assumption 2.2 hold, and let x_\top^* and y_2^* be defined as in Assumption 2.3 and (3.5), respectively. Then, y_2^* satisfies

$$\frac{\|x_\top^*\|_2^2}{\|y_2^*\|_2^2} \sim \chi_{d-d_e+1}^2.$$

Proof. The result almost immediately follows from Corollary 3.6 and Lemma A.15. These yield

$$\frac{\|z^*\|^2}{\|y_2^*\|^2} \sim \chi_{d-d_e+1}^2.$$

The result is implied by $\|z^*\| = \|x_\top^*\|$ (see Remark 3.2). \square

The above result is equivalent to saying that $\|y_2^*\|^2/\|x_\top^*\|^2$ follows the inverse chi-squared distribution with $d - d_e + 1$ degrees of freedom (see Definition A.7). The theorem reveals a linear dependence of $\|y_2^*\|$ on $\|x_\top^*\|$; larger values of $\|x_\top^*\|$ contribute to the increase in the likelihood of y_2^* being further away from the origin. The theorem also suggests that $\|y_2^*\|$ is independent of D as long as $\|x_\top^*\|$ is fixed.

COROLLARY 3.8 Let Assumption 2.2 hold. Let x_\top^* and y_2^* be defined as in Assumption 2.3 and (3.5), respectively. Then,

$$\mathbb{P}[\|y_2^*\|_2 \leq \delta] = \mathbb{P}\left[\chi_{d-d_e+1}^2 \geq \frac{\|x_\top^*\|_2^2}{\delta^2}\right]$$

for any $\delta > 0$.

Proof. For any $\epsilon > 0$, we have

$$\mathbb{P}\left[\|y_2^*\|_2 \leq \frac{\|x_\top^*\|_2}{\epsilon}\right] = \mathbb{P}\left[\frac{\|x_\top^*\|_2^2}{\|y_2^*\|_2^2} \geq \epsilon^2\right] = \mathbb{P}[\chi_{d-d_e+1}^2 \geq \epsilon^2],$$

where the second equality follows from Theorem 3.7. By letting $\epsilon = \|x_\top^*\|_2/\delta$, we obtain the result. \square

COROLLARY 3.9 Let Assumption 2.2 hold, and let \mathbf{x}_\top^* and \mathbf{y}_2^* be defined as in Assumption 2.3 and (3.5), respectively. Provided that $d - d_e > 1$, we have

$$\mathbb{E}[\|\mathbf{y}_2^*\|^2] = \frac{\|\mathbf{x}_\top^*\|^2}{d - d_e - 1}. \quad (3.8)$$

Proof. Let W follow the inverse chi-squared distribution with $d - d_e + 1$ degrees of freedom. Then, $W \stackrel{\text{law}}{=} \|\mathbf{y}_2^*\|^2 / \|\mathbf{x}_\top^*\|^2$ by Theorem 3.7. By applying Lemma A.8, we obtain

$$\mathbb{E}[\|\mathbf{y}_2^*\|^2] = \mathbb{E}[\|\mathbf{x}_\top^*\|^2 W] = \frac{\|\mathbf{x}_\top^*\|^2}{d - d_e - 1}$$

for $d - d_e > 1$. □

The expected value in (3.8) is inversely proportional to $d - d_e$. In other words, for a fixed d_e , larger values of the dimension of the embedding subspace bring \mathbf{y}_2^* closer to the origin. This observation indicates that the increase in d allows us to decrease δ while the probability of $\mathbf{y}_2^* \in \mathcal{Y}$ is kept constant.

3.2.2 The probability density function The following theorem derives the probability density function of \mathbf{y}_2^* .

THEOREM 3.10 Let Assumption 2.2 hold, and let \mathbf{x}_\top^* and \mathbf{y}_2^* be defined as in Assumption 2.3 and (3.5), respectively. Then, the probability density function of \mathbf{y}_2^* is given by

$$g^*(\mathbf{y}) = \pi^{-d/2} \left(\frac{\Gamma(d/2)}{\Gamma(n/2)} \right) \left(\frac{\|\mathbf{x}_\top^*\|}{\sqrt{2}} \right)^n (\mathbf{y}^T \mathbf{y})^{-(n+d)/2} e^{-\|\mathbf{x}_\top^*\|^2 / (2\mathbf{y}^T \mathbf{y})},$$

where $n = d - d_e + 1$.

Proof. Corollary 3.6 and Lemma A.17 imply that the p.d.f. of \mathbf{y}_2^* is given by

$$g^*(\mathbf{y}) = \pi^{-d/2} \left(\frac{\Gamma(d/2)}{\Gamma(n/2)} \right) \left(\frac{\|\mathbf{z}^*\|}{\sqrt{2}} \right)^n (\mathbf{y}^T \mathbf{y})^{-(n+d)/2} e^{-\|\mathbf{z}^*\|^2 / (2\mathbf{y}^T \mathbf{y})}.$$

By using the equation $\|\mathbf{z}^*\| = \|\mathbf{x}_\top^*\|$, we obtain the desired result. □

Figure 3 illustrates the p.d.f. of two-dimensional \mathbf{y}_2^* . The shape of the p.d.f. resembles a volcano with the mass concentrated at a certain distance from the origin suggesting that \mathbf{y}_2^* is unlikely to be neither too close to nor too distant from the origin. We also note that the p.d.f. is independent of D .

4. Bounding the success of the reduced problem

This section is the culmination of this paper's analysis. Based on the results established earlier, we derive a bound for the probability of success of (RP).

The following theorem presents a notable connection between the success of (RP) and the chi-squared distribution.

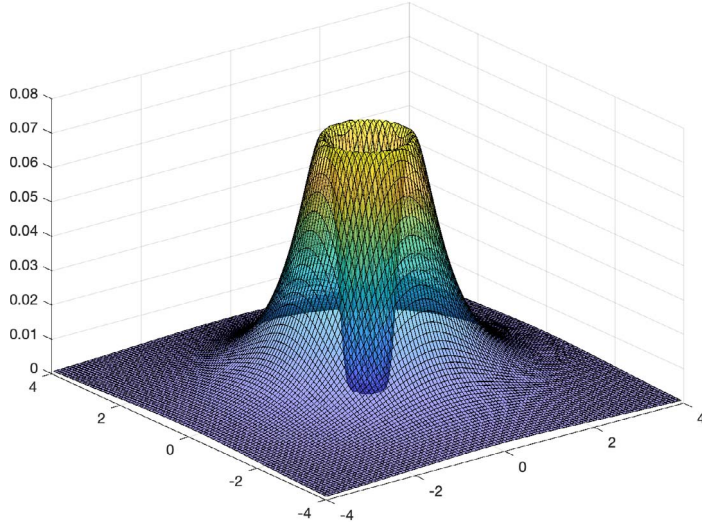


FIG. 3. The illustration of the p.d.f. of y_2^* for $d = 2$, $k = 2$ and $\mathbf{x}_\top^* = [1 \ 1]^T$.

THEOREM 4.1 Let Assumption 2.2 hold. Then, for any $\delta > 0$, we have

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq \mathbb{P}\left[\chi_{d-d_e+1}^2 \geq \frac{\min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{x}^*\|_2^2}{\delta^2}\right]. \quad (4.1)$$

Proof. Note the following relationship between the probabilities:

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq \mathbb{P}[\mathbf{y}_2^* \in \mathcal{Y}] = \mathbb{P}[\|\mathbf{y}_2^*\|_\infty \leq \delta] \geq \mathbb{P}[\|\mathbf{y}_2^*\|_2 \leq \delta], \quad (4.2)$$

where the first inequality follows from Lemma 3.5 and where the second inequality is implied by $\|\mathbf{y}_2^*\|_\infty \leq \|\mathbf{y}_2^*\|_2$. By applying Corollary 3.8 to the last probability in (4.2) and using the definition of \mathbf{x}_\top^* given in Assumption 2.3, we obtain

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq \mathbb{P}\left[\chi_{d-d_e+1}^2 \geq \frac{\|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2^2}{\delta^2}\right] \quad (4.3)$$

for any $\delta > 0$ and any $\mathbf{x}^* \in \mathcal{G}$ such that $\|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2 \neq 0$. Note that (4.3) also holds for $\|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2 = 0$ since, in this case, (RP) is successful with probability 1 (see the discussion preceding Assumption 2.3). Hence, (4.3) holds for any $\mathbf{x}^* \in \mathcal{G}$, which then implies

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq \max_{\mathbf{x}^* \in \mathcal{G}} \mathbb{P}\left[\chi_{d-d_e+1}^2 \geq \frac{\|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2^2}{\delta^2}\right] = \mathbb{P}\left[\chi_{d-d_e+1}^2 \geq \frac{\min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2^2}{\delta^2}\right],$$

where the equality follows from the fact that the tail distribution $\mathbb{P}[X > x]$ of any random variable X is a monotonically decreasing function in x .

In what follows, we show that $\min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2^2 = \min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{x}^*\|_2^2$. Define sets $\mathcal{Z} = \{\mathbf{z} \in \mathbb{R}^d : \mathbf{U}\mathbf{z} = \mathbf{U}\mathbf{U}^T \mathbf{x}^*, \mathbf{x}^* \in \mathcal{G}\}$ and $\mathcal{S} = \{\mathbf{U}\mathbf{z} + \mathbf{V}\mathbf{c} : \mathbf{z} \in \mathcal{Z}, \mathbf{c} \in \mathbb{R}^{D-d_e}\}$, where \mathbf{V} is defined in Assumption 2.3. First, we establish that $\mathcal{G} = \mathcal{S}$ by showing that $\mathcal{G} \subseteq \mathcal{S}$ and that $\mathcal{S} \subseteq \mathcal{G}$.

Let $\mathbf{x}^* \in \mathcal{G}$. We can write $\mathbf{x}^* = \mathbf{U}\mathbf{U}^T \mathbf{x}^* + \mathbf{V}\mathbf{V}^T \mathbf{x}^*$ since $\mathbf{U}\mathbf{U}^T + \mathbf{V}\mathbf{V}^T = \mathbf{I}$. Let $\mathbf{z} = \mathbf{U}^T \mathbf{x}^*$ and $\mathbf{c} = \mathbf{V}^T \mathbf{x}^*$ and note that $\mathbf{z} \in \mathcal{Z}$ and $\mathbf{c} \in \mathbb{R}^{D-d_e}$. Hence, $\mathbf{x}^* \in \mathcal{S}$, which proves that $\mathcal{G} \subseteq \mathcal{S}$.

Let $\mathbf{x}^* \in \mathcal{S}$. Then, $\mathbf{x}^* = \mathbf{U}\mathbf{z} + \mathbf{V}\mathbf{c}$ for some $\mathbf{z} \in \mathcal{Z}$ and $\mathbf{c} \in \mathbb{R}^{D-d_e}$. We have

$$f(\mathbf{x}^*) = f(\mathbf{U}\mathbf{z} + \mathbf{V}\mathbf{c}) = f(\mathbf{U}\mathbf{z}) = f(\mathbf{U}\mathbf{U}^T \mathbf{x}^*) = f(\mathbf{x}_\top^*) = f^*,$$

where the second equality follows from the assumption that f has low effective dimensionality and the fact that $\mathbf{V}\mathbf{c} \in \mathcal{T}^\perp$, the fourth equality follows from the definition of \mathbf{x}_\top^* (given in Assumption 2.3) and the last equality follows from (2.2). Hence, by definition of \mathcal{G} , $\mathbf{x}^* \in \mathcal{G}$. This proves that $\mathcal{S} \subseteq \mathcal{G}$.

Finally, we have

$$\begin{aligned} \min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{x}^*\|_2^2 &= \min_{\mathbf{x}^* \in \mathcal{S}} \|\mathbf{x}^*\|_2^2 = \min_{\mathbf{z} \in \mathcal{Z}, \mathbf{c} \in \mathbb{R}^{D-d_e}} \|\mathbf{U}\mathbf{z} + \mathbf{V}\mathbf{c}\|_2^2 && \text{(since } \mathcal{G} = \mathcal{S} \text{ and by definition of } \mathcal{S}\text{)} \\ &= \min_{\mathbf{z} \in \mathcal{Z}, \mathbf{c} \in \mathbb{R}^{D-d_e}} \|\mathbf{U}\mathbf{z}\|_2^2 + \|\mathbf{V}\mathbf{c}\|_2^2 && \text{(since } \mathbf{U}^T \mathbf{V} = \mathbf{V}^T \mathbf{U} = \mathbf{0}\text{)} \\ &= \min_{\mathbf{z} \in \mathcal{Z}, \mathbf{c} \in \mathbb{R}^{D-d_e}} \|\mathbf{U}\mathbf{z}\|_2^2 + \|\mathbf{c}\|_2^2 && \text{(since } \mathbf{V} \text{ is orthogonal)} \\ &= \min_{\mathbf{z} \in \mathcal{Z}} \|\mathbf{U}\mathbf{z}\|_2^2 + \min_{\mathbf{c} \in \mathbb{R}^{D-d_e}} \|\mathbf{c}\|_2^2 \\ &= \min_{\mathbf{z} \in \mathcal{Z}} \|\mathbf{U}\mathbf{z}\|_2^2 + 0 \\ &= \min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{U}\mathbf{U}^T \mathbf{x}^*\|_2^2 && \text{(by definition of } \mathcal{Z}\text{)}. \end{aligned}$$

□

Using Theorem 4.1, one can now bound the success of (RP) by applying any tail bound on the chi-squared distribution. We use the bound derived in Lemma A.6.

COROLLARY 4.2 Let Assumption 2.2 hold, and let $\mu = \min_{\mathbf{x}^* \in \mathcal{G}} \|\mathbf{x}^*\|_2$. Then, for any $\delta > 0$, we have

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq 1 - C(n) \left(1 + \frac{n}{2} e^{-\mu^2/(2\delta^2)}\right) \left(\frac{\mu}{\sqrt{2}\delta}\right)^n, \quad (4.4)$$

where $n = d - d_e + 1$ and

$$C(n) = \frac{4}{n(n+2)\Gamma(n/2)}.$$

Proof. Lemma A.6 implies that

$$\mathbb{P}[\chi_n^2 \geq \epsilon^2] \geq 1 - C(n) \left(1 + \frac{n}{2} e^{-\epsilon^2/2}\right) (\epsilon^2/2)^{n/2} \quad (4.5)$$

for any $\epsilon > 0$. By letting $\epsilon = \mu/\delta$ and applying (4.5) to (4.1), we obtain the wished bound. \square

Let R^* denote the right-hand side of (4.4). First, we note that R^* is a function of μ/δ and $d - d_e$. The bound reveals a linear relationship between μ and δ ; scaling μ and δ by the same factor does not affect the value of R^* . Furthermore, observe that for smaller values of μ or larger values of δ , R^* is closer to 1. Numerical experiments show that for large values of n and/or μ/δ , the bound (4.4) is less tight; this is also signified by the asymptotic behaviour of R^* , $R^* \rightarrow -\infty$ monotonically as $\mu/\delta \rightarrow \infty$ making the bound useless for large enough μ/δ .

It is remarkable that R^* has no dependence on D , the dimension of the original optimization problem. This implies that larger D does not diminish the success of the reduced problem as long as μ and d_e are unchanged. Dependence of R^* on $d - d_e$ indicates that the success is determined by the value of d relative to d_e and not so much by the individual values of d and d_e . Larger (smaller) values of d with respect to d_e require smaller (larger) δ if R^* is kept constant; knowing this fact is crucial when initializing values of d and δ in practice. It displays a convenient interplay between d and δ allowing more flexibility in choosing one vs. another.

Previous bounds. One can derive similar bounds for the success of (RP) by bounding $\mathbb{P}[\|\mathbf{y}_2^*\| \leq \delta]$ in (4.2) using the Cauchy–Schwarz inequality. Since $\mathbf{y}_2^* = \mathbf{B}^T(\mathbf{B}\mathbf{B}^T)^{-1}\mathbf{z}^*$, we have

$$\|\mathbf{y}_2^*\| \leq \|\mathbf{B}^T(\mathbf{B}\mathbf{B}^T)^{-1}\| \cdot \|\mathbf{z}^*\|.$$

By using the fact that $\|\mathbf{B}^T(\mathbf{B}\mathbf{B}^T)^{-1}\| = 1/s_{\min}(\mathbf{B}^T)$, where $s_{\min}(\mathbf{B}^T)$ denotes the smallest singular value of \mathbf{B}^T , we obtain

$$\mathbb{P}[\|\mathbf{y}_2^*\| \leq \delta] \geq \mathbb{P}\left[\frac{\|\mathbf{z}^*\|}{s_{\min}(\mathbf{B}^T)} \leq \delta\right].$$

We can now use any suitable tail bound for the smallest singular value of the Gaussian matrix to bound the latter probability.

Wang *et al.* [49], by applying the above technique and the result in [14] to bound the singular value, derived the following bound:

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq 1 - \frac{\mu\sqrt{d_e}}{\delta}.$$

Their derivation is predicated on the assumptions that $d = d_e$ and that \mathcal{T} is spanned by the standard basis vectors. Sanyang & Kabán [43] extended Wang *et al.*'s bound to any δ satisfying $\delta > \|\mathbf{x}_1^*\|/(\sqrt{d} - \sqrt{d_e})$. Using the bound in [10] for $s_{\min}(\mathbf{B}^T)$, they showed that

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq 1 - e^{-(\sqrt{d} - \sqrt{d_e} - \mu/\delta)^2/2}.$$

One can also use Rudelson & Vershynin's [41, Theorem 1.1] bound to obtain

$$\mathbb{P}[(\text{RP}) \text{ is successful}] \geq 1 - \left(\frac{C\mu}{\delta(\sqrt{d} - \sqrt{d_e} - 1)}\right)^{d-d_e+1} - e^{-cd},$$

where $C, c > 0$ are absolute constants. This bound shows dependence of the probability on the difference $d - d_e$, which also manifests in our bound. Rudelson and Vershynin's bound cannot be used for practical purposes due to the unknown C and c ; we require explicit bounds to define the size of \mathcal{Y} .

Unlike the bounds of Wang *et al.* [49] and Sanyang & Kaban [43], Corollary 3.8 is applicable to any $d \geq d_e$ and an arbitrary subspace \mathcal{T} . Moreover, using the exact distribution of $\|y_2^*\|$ given in Corollary 3.8, we circumvent the application of the intermediate Cauchy–Schwarz and bound the distribution of $\|y_2^*\|$ directly.

Affine random embeddings. It is not difficult to extend (RP) to affine random subspace embeddings. In the affine case, we replace x by $Ay + p$, where $p \in \mathbb{R}^D$ is a fixed point. The reduced optimization problem is then given by

$$\begin{aligned} \min \quad & f(Ay + p) \\ \text{subject to} \quad & y \in \mathcal{Y}. \end{aligned}$$

The results that apply to the linear embeddings also apply to the affine embeddings after minor adjustments. Theorem 2.4, for example, can be easily extended to the affine case to show that the intersection between $p + \text{range}(A)$ and \mathcal{G} takes place with probability 1 if $d \geq d_e$. The affine version of the results are provable with the same assumptions except for a minor alteration in Assumption 2.3: the condition $x_\top^* \neq 0$ changes to $x_\top^* \neq p$. To obtain the affine versions of Theorems 3.7 and 3.10, replace x_\top^* with $x_\top^* - p_\top$, where $p_\top = UU^T p$ is the orthogonal projection of p onto \mathcal{T} . For the affine versions of Theorem 4.1 and Corollary 4.2, replace $\min_{x^* \in \mathcal{G}} \|x^*\|$ with $\min_{x^* \in \mathcal{G}} \|x^* - p\|$.

5. Numerical experiments

5.1 Choices of (RP) parameters

The present section aims to test numerically the quality of the bound (4.4). We will also use the results of this section to select suitable pairs of parameters d and δ for (RP) in the numerical experiments later.

Suppose that we are given a function f satisfying Assumption 2.2 with the set of global minimizers \mathcal{G} consisting of only one connected component. Let x_\top^* for f be defined as in Assumption 2.3 and z be defined by the equation $Uz = x_\top^*$. We also define $\mu := \min_{x^* \in \mathcal{G}} \|x^*\|$ and note that $\mu = \|x_\top^*\|$.

We test (4.4) for f by contrasting the left-hand side of (4.4) (denoted by L^*) to its right-hand side (denoted by R^*). We compare L^* and R^* for four different values of $d - d_e$, namely 0, 1, 2 and 3. For each value of $d - d_e$, we express R^* as a function of $\bar{\delta} := \delta/\mu$ and using its closed form we plot R^* for $\bar{\delta} \in [0.02, 10]$. We do not have a closed form expression for L^* , but we can approximate it numerically. In what follows, we describe how this could be done. We start by writing

$$\begin{aligned} L^* := \mathbb{P}[(\text{RP}) \text{ is successful}] &= \mathbb{P}[\exists y \in [-\delta, \delta]^d : Ay \in \mathcal{G}] = \mathbb{P}[\exists y \in [-\delta, \delta]^d : \bar{B}y = z] \\ &= \mathbb{P}[\exists y \in [-\bar{\delta}, \bar{\delta}]^d : \bar{B}y = \bar{z}], \end{aligned} \quad (5.1)$$

where \bar{B} denotes a $d_e \times d$ Gaussian matrix and $\bar{z} = z/\mu$. Here, the second equality follows from Definition 1.2 and the third equality follows from Theorem 3.1 and the fact that \mathcal{G} has only one connected component. Note that $\|\bar{z}\| = 1$ since $\|z\| = \|x_\top^*\| = \mu$ (see Remark 3.2). We assign \bar{z} to a

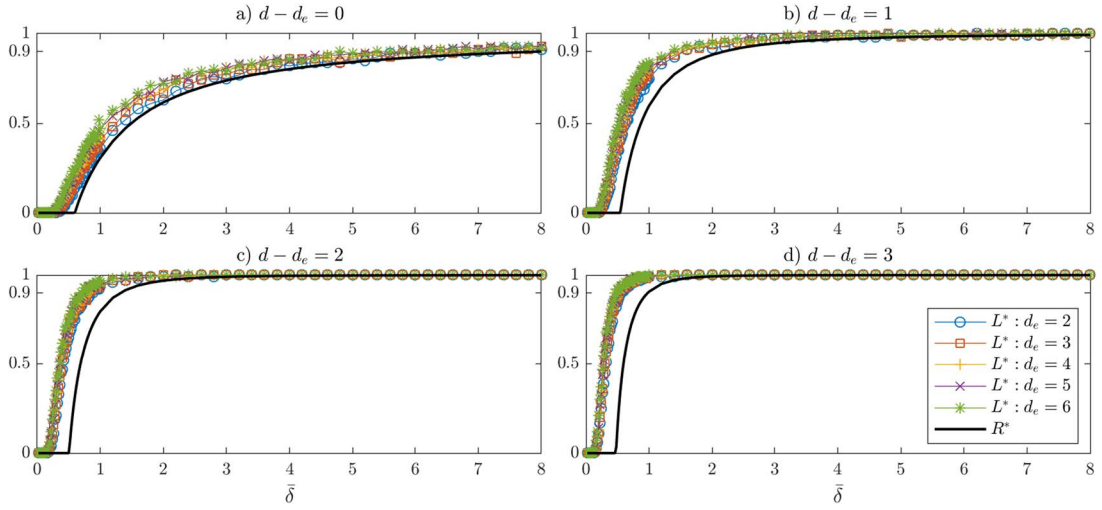


FIG. 4. The four plots depict the function $R^*(\bar{\delta})$ and the estimates of $L^*(\bar{\delta})$; each plot corresponds to a particular value of $d - d_e \in \{0, 1, 2, 3\}$. Each plot contains estimates of $L^*(\bar{\delta})$ for $d_e = 2, 3, 4, 5, 6$.

random vector with unit norm and keep \bar{z} fixed throughout the experiment⁷. For each $\bar{\delta} \in [0.02, 10]$, we generate a thousand Gaussian matrices \bar{B} and estimate the latter probability in (5.1) as the proportion of instances for which the statement under the probability is true. Unlike for R^* , L^* depends on individual values of d and d_e . We plot the estimates for L^* for the following values of $d - d_e$: 0, 1, 2, 3 and, in each plot, we repeat the experiment for $d_e = 2, 3, 4, 5, 6$. The plots are presented in Fig. 4.

Numerical findings. The plots in Fig. 4—confirming the conclusions of Corollary 4.2—suggest that the variation in success of (RP) is mainly determined by the value of $d - d_e$; the larger is the difference, the higher is the probability of success of (RP) for a given $\bar{\delta}$. These curves, being independent of μ , can be used to find suitable $\bar{\delta}$ for any problem for the corresponding values of $d - d_e$; the size of the \mathcal{Y} box, δ , can then be set to $M\bar{\delta}$ if an upper bound M on μ is known.

Choosing d , μ and δ in practice. When it comes to the numerical application of (RP) in practice, initialization of parameters d and δ might be problematic. From the theoretical discussions above, we learned that the parameters d and δ must be defined based on d_e and μ , the values of which are typically unknown in practice, for example, for black-box functions. We circumvent this issue by estimating d_e and μ rather than trying to calculate their exact values; note that all we need is an upper bound d on d_e . The parameter d_e or an upper bound may be known from prior studies or can be found with active subspace identification methods (see, e.g. [9]); these use gradients of f to estimate d_e . Another approach that can avoid the need to use gradients of active subspace methods is the use of some numerical trial-and-error procedure, which, for example, may involve a systematic increase of d until no significant changes in the objective function are observed [43].

⁷ Note that the results of the experiment are invariant of the choice of \bar{z} as long as its norm is fixed. Let z_1 and z_2 be two fixed vectors with unit norm. Consider two systems: $By = z_1$ and $By = z_2$. Note that z_2 can be written as Qz_1 for some orthogonal $Q \in \mathbb{R}^{d_e \times d_e}$. Then, the second system becomes $Q^T B y = z_1$ and this generates vectors y with the same distribution as the first system since $Q^T B$ is also Gaussian.

Estimating μ can be a harder task. A rough estimate for μ can be obtained if the search in the original space is restricted to a certain domain; a trivial upper bound in this case is given by the maximum distance between the origin and the boundary of the domain. The search domain that is commonly imposed to practically solve unconstrained optimization problems is box constraints, such as $\mathcal{X} = [-1, 1]^D$ for which $\mu \leq \sqrt{D}$. In Appendix C, we test REGO assuming that \sqrt{D} is the best bound known for μ . To compensate for unknown μ and depending on available computational budget, one could also try increasing δ or d gradually to explore larger regions in \mathbb{R}^d (and correspondingly in \mathbb{R}^D). Further discussions of these issues are given in the Conclusions.

ALGORITHM 5.1 (Random Embeddings for Global Optimization (REGO applied to (P))).

1. Initialise d and δ and define $\mathcal{Y} = [-\delta, \delta]^d$
2. Generate a $D \times d$ Gaussian matrix A
3. Apply a global optimization solver (e.g. BARON, DIRECT, KNITRO) to (RP) until a termination criterion is satisfied, and define y_{\min} to be the generated (approximate) solution of (RP).
4. Reconstruct $x_{\min} = Ay_{\min}$

5.2 Testing REGO with state-of-the-art global solvers

Algorithms.

The algorithm for the random embeddings method named REGO is outlined in Algorithm 5.1. Below, we give the descriptions of the three state-of-the-art solvers we use to test REGO.

DIRECT [18, 21, 27] version 4.0 (Dividing RECTangles) is a deterministic⁸ global optimization solver first introduced in [27] as an extension of Lipschitzian optimization. DIRECT does not require information about the gradient or about the Lipschitz constant and, hence, can be used for black-box functions. DIRECT divides the search domain into rectangles and evaluates the function at the centre of each rectangle. Based on the previously sampled points, DIRECT carefully decides what rectangle to divide next balancing between local and global searches. Jones *et al.* [27] showed that DIRECT is guaranteed to converge to global minimum, but convergence may sometimes be slow.

BARON [42, 45] version 17.10.10 (Branch-And-Reduce Optimization Navigator) is a branch- and-bound type global optimization solver for nonlinear and mixed-integer nonlinear programs. To provide lower and upper bounds for each branch, BARON utilizes algebraic structure of the objective function. It also includes a preprocessing step where it performs a multi-start local search to obtain a tight global upper bound. In comparison to other existing global solvers, BARON was demonstrated to be the most robust and fastest [35]. However, BARON accepts only a few (general) classes of functions⁹ including polynomial, exponential, logarithmic, etc., and, unlike DIRECT, it is unable to optimize black-box functions.

KNITRO [6] version 10.3.0 is a large-scale nonlinear local optimization solver capable of handling problems with hundreds of thousands of variables. KNITRO allows to solve problems using one of the four algorithms: two interior point type methods (direct and conjugate gradient) and two active set type methods (active set and sequential quadratic programming). In contrast to BARON and DIRECT, which specialize on finding global minima, KNITRO focuses on finding local solutions. Nonetheless,

⁸ Here, we refer to the predictable behaviour of the solver given a fixed set of parameters.

⁹ For instance, BARON cannot be applied to problems which include trigonometric functions.

TABLE 1 The table outlines the experimental setup for the three solvers. In the table, f is a function with low effective dimensionality d_e and the global minimum f^* and ϵ is set to 10^{-3}

| | DIRECT | BARON | KNITRO |
|---------------------------------------|---|--|---|
| Measure of computational cost | Function evaluations | CPU seconds | Function evaluations, CPU seconds |
| Budget per problem | $10000 \times d_e$ function evaluations | $200 \times d_e$ CPU seconds | $20 \times d_e$ starting points |
| Convergence criteria (see Remark 5.2) | $f_D^* \leq f^* + \epsilon$ | Convergence: $f_B^U \leq f^* + \epsilon$ Convergence _{opt} : $f_B^U \leq f_B^L + \epsilon$ | $f_K^* \leq f^* + \epsilon$ |
| Termination criteria | Either on budget or if \mathbf{x}_D^* satisfies the convergence criteria | Either on budget or if f_B^U and f_B^L satisfy the convergence _{opt} criteria | On budget |
| Additional options | <code>options.testflag=1</code> <code>options.maxits=Inf</code> <code>options.globalmin=f*</code> | <code>npsol = 9 numloc = 0</code> <code>BrVarStra = 1</code> <code>BrPtStra = 1</code> | Default options. Derivatives allowed. Use of multi-start through <code>ms_enable=1</code> . |

KNITRO has multi-start capabilities, i.e. it solves a problem locally multiple times every time starting from a different point in the feasible domain. It is this feature that we make use of in the experiments.

Generating the test set. Our test set of functions with low effective dimensionality will be derived from 19 global optimization problems (of dimensions 2–6) with known global minima [16, 23, 44], some of which are from the Dixon–Szego set [12]. The list of the problems is given in Table B3, Appendix B.

Below, we describe the method adopted from Wang *et al.* [49] to generate high-dimensional functions with low effective dimensionality. Let $\bar{g}(\bar{\mathbf{x}})$ be any function from Table B3; let d_e be its dimension, and let the given domain be scaled to $[-1, 1]^{d_e}$. We create a D -dimensional function $g(\mathbf{x})$ by adding $D - d_e$ fake dimensions to $\bar{g}(\bar{\mathbf{x}})$, $g(\mathbf{x}) = \bar{g}(\bar{\mathbf{x}}) + 0 \cdot x_{d_e+1} + 0 \cdot x_{d_e+2} + \dots + 0 \cdot x_D$. We further rotate the function by applying a random orthogonal matrix \mathbf{Q} to \mathbf{x} to obtain a non-trivial constant subspace. The final form of the function we test is given as

$$f(\mathbf{x}) = g(\mathbf{Q}\mathbf{x}). \quad (5.2)$$

Note that the first d_e rows of \mathbf{Q} now span the effective subspace \mathcal{T} of $f(\mathbf{x})$. Furthermore,

$$\mu := \min_{\mathbf{x} \in \mathcal{G}_1} f(\mathbf{x}) = \min_{\bar{\mathbf{x}} \in \mathcal{G}_2} \bar{g}(\bar{\mathbf{x}}) \leq \sqrt{d_e}, \quad (5.3)$$

where \mathcal{G}_1 and \mathcal{G}_2 are the sets of global minimizers of f and \bar{g} , respectively.

For each problem in the test set, we generate three functions f as defined in (5.2) one for each $D = 10, 100, 1000$. We will tackle (P) for each f both directly (we call it *no embedding*) and applying REGO outlined in Algorithm 5.1.

Experimental setup (REGO). We compare ‘no embedding’ and REGO using the three solvers above. Let g_i, s_j, n_j and D_k denote the i th function in the problem set ($g_1 = \text{Beale}$, etc.; see Table B3), j th solver ($s_1 = \text{DIRECT}$, $s_2 = \text{BARON}$, $s_3 = \text{KNITRO}$), the total number of problems in the problem set solvable by j th solver ($n_1 = 19, n_2 = 15, n_3 = 18$) and k th ambient dimension ($D_1 = 10, D_2 = 100, D_3 = 1000$), respectively. Let f_{ik} denote the D_k -dimensional function with low effective dimensionality constructed from g_i as described previously.

Within ‘no embedding’ framework, for each pair (s_j, D_k) , we solve f_{ik} for $i = 1, 2, \dots, n_j$ with solver s_j and record the proportion of the problems that attain convergence (see definition in Table 1).

For each f_{ik} ($1 \leq i \leq n_j, 1 \leq k \leq 3$), we apply REGO 100 times every time with a different Gaussian matrix. Thus, in total, for each pair (s_j, D_k) we solve $n_j \times 100$ problems. We record the proportion of problems that attain convergence (Table 1) out of these $n_j \times 100$ problems.

We also record the number of function evaluations (for DIRECT and KNITRO) and CPU time (for all the three solvers) spent before termination within the two frameworks. For each (s_j, D_k) , function evaluations and time are averaged out over $n_j \times 100$ problems within REGO and over n_j problems within ‘no embedding’.

We conduct the above experiment for REGO with the following pairs of parameters (d, δ) : $(d_e, 8.0 \times \sqrt{d_e}), (d_e + 1, 2.2 \times \sqrt{d_e}), (d_e + 2, 1.3 \times \sqrt{d_e})$ and $(d_e + 3, 1.0 \times \sqrt{d_e})$. Here, each δ was set to $M\bar{\delta}$, where $M = \sqrt{d_e}$ is an upper bound on μ (see (5.3)) and the value for $\bar{\delta}$ was chosen as the smallest $\bar{\delta}$ that gives at least 90% chance of success based on the curve of R^* in Fig. 4.

Experimental setup (solvers). Due to the difference in algorithmic procedures of the solvers, they allow different budget constraints and have different convergence and termination criteria; we present these in Table 1.

REMARK 5.2 DIRECT, at its every iteration, stores f_D^* —the minimum value of f so far found. BARON, at its every iteration, stores f_B^U and f_B^L —smallest upper bound and largest lower bound so far found for f . As for KNITRO, $f_K^* = \min\{f(\mathbf{A}\mathbf{y}_1^*), f(\mathbf{A}\mathbf{y}_2^*), \dots, f(\mathbf{A}\mathbf{y}_l^*)\}$, where l is the number of starting points and where $\{\mathbf{y}_i^*\}_{1 \leq i \leq l}$ are the local solutions produced by the multi-start procedure.

REMARK 5.3 The experiments are done not to compare solvers but to contrast ‘no embedding’ with REGO. All the experiments were run in MATLAB on the 16 cores (2×8 Intel with hyper-threading) Linux machines with 256GB RAM and 3300 MHz speed.

5.3 Numerical results

(RP) successful. We record the proportion of instances for which (RP) is successful. Table 2 presents these percentages for each particular choice of d and D averaged over 19 problems in the test set. We observe that the percentages are very high and appear to be independent of D supporting the conclusions of Corollary 4.2.

REGO vs. no embedding. The results of the experiment comparing REGO and ‘no embedding’ are presented in Figs 5–7 for DIRECT, BARON and KNITRO, respectively. These figures compare average proportions of converged solutions and computational costs produced by REGO and ‘no embedding’ frameworks for $D = 10, 100, 1000$.

TABLE 2 The table shows average percentages of problems for which (RP) is successful

| $d \setminus D$ | 10 | 100 | 1000 |
|-----------------|------|------|------|
| $d_e + 0$ | 97.2 | 97.8 | 97.3 |
| $d_e + 1$ | 99.1 | 98.9 | 99.3 |
| $d_e + 2$ | 99.5 | 99.6 | 99.8 |
| $d_e + 3$ | 100 | 99.9 | 99.8 |

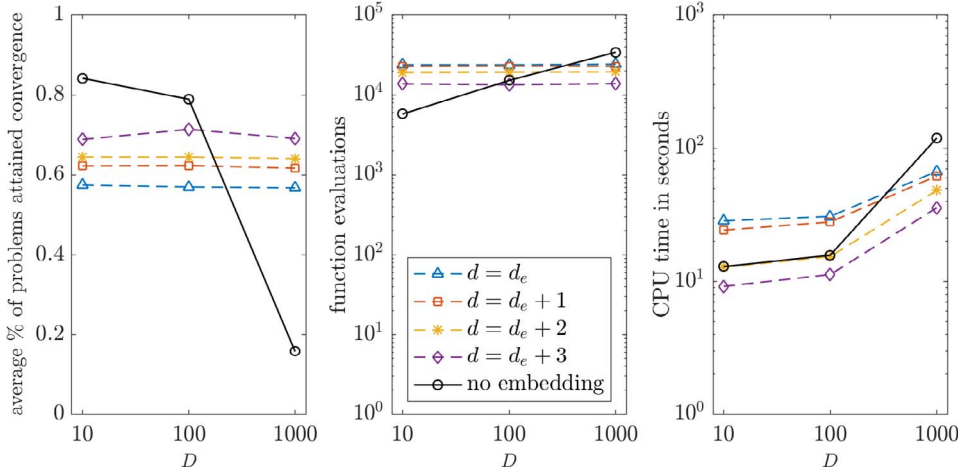


FIG. 5. REGO vs. 'no embedding' with DIRECT: comparison of frequency of convergence, log of average function evaluations and log of average CPU time (in seconds).

DIRECT (Fig. 5). For all the four initializations of REGO, we observe that the average proportions of problems that attained convergence (see definition in Table 1) are invariant with respect to the ambient dimension. This frequency of convergence is higher within 'no embedding' for $D = 10, 100$, but exhibits a significant drop for $D = 1000$. The average function evaluation count is maintained within REGO but doubles within 'no embedding' for a tenfold increase in D . Growth in CPU time takes place within both frameworks, being the highest for 'no embedding'.

BARON (Fig. 6). In comparison with 'no embedding', the frequency of convergence $_{opt}$ is higher within REGO in most cases. We note that BARON's both convergence and convergence $_{opt}$ exhibit invariance with respect to D within REGO. As for 'no embedding', we observe a decrease in the frequencies of both convergence and convergence $_{opt}$. In addition, we observe an increase in CPU time spent within 'no embedding', while the time is almost constant within REGO.

KNITRO (Fig. 7). We see that the proportion of solved problems is invariant with respect to the ambient dimension within REGO and, surprisingly, within 'no embedding' as well. However, the average number of function evaluations and time spent differ significantly between the two frameworks. With REGO, the average number of function evaluations remain at the same level for all D . Average time grows within both frameworks but at a higher rate for 'no embedding'. The average time differs by a factor of 70 for $D = 1000$ in favour of REGO. We think that the growth in time within REGO is due to more costly function and derivative evaluations for larger D .

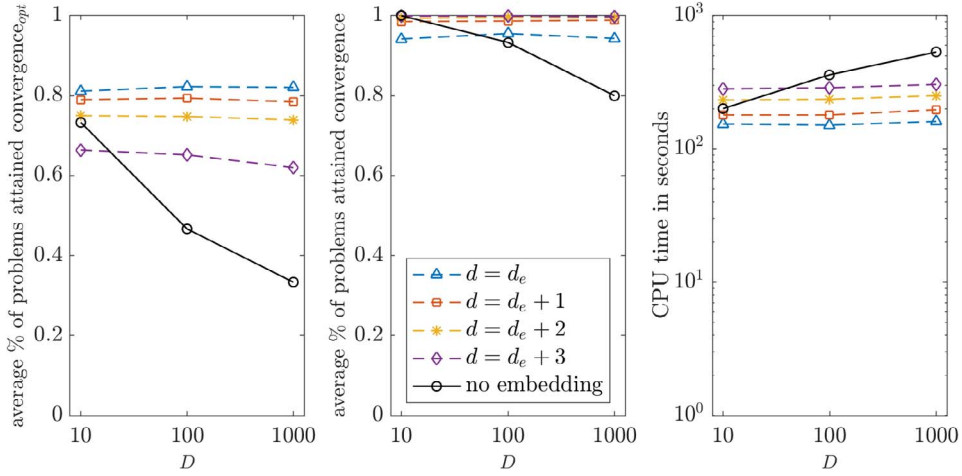


FIG. 6. REGO vs. 'no embedding' with BARON: comparison of frequency of convergence_{opt}/convergence and average CPU time (in seconds).

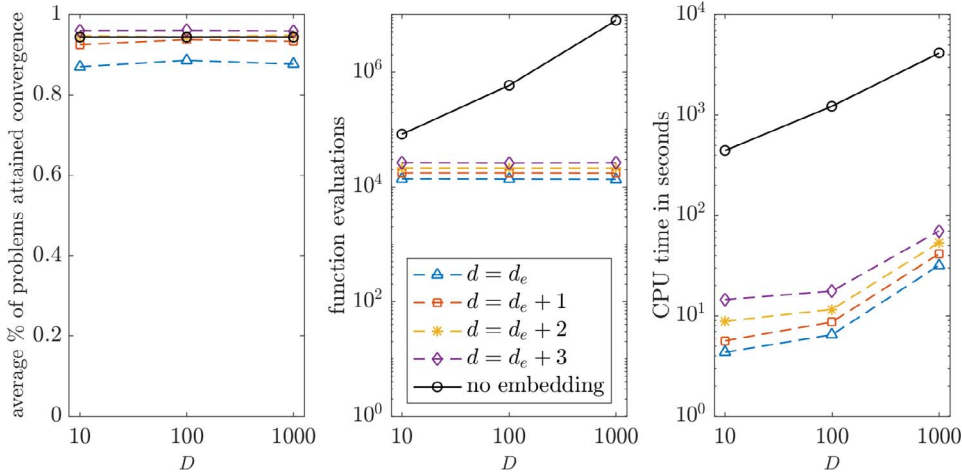


FIG. 7. REGO vs. 'no embedding' with KNITRO: comparison of frequency of convergence, log of average function evaluations and log of average CPU time (in seconds).

5.4 Summary of numerical results

1. (Effects of parameter choices) Our experiments clearly show that the choice of d and δ has a considerable effect on convergence and computational cost of REGO and that good choices of (d, δ) are dependent on the underlying solver. For example, BARON achieved the highest proportion of convergence_{opt} within least amount of time for $(d, \delta) = (d_e, 8\sqrt{d_e})$, whereas DIRECT performed best for $(d_e + 3, \sqrt{d_e})$. KNITRO produced the highest proportion of convergence and the worst time for $(d_e + 3, \sqrt{d_e})$ and the lowest proportion of convergence and the best time for $(d_e, 8\sqrt{d_e})$.

2. (Scalability) Within REGO, the proportion of problems solved and/or the number of function evaluations are generally invariant with respect to the ambient dimension D . REGO displays good scalability for all three solvers.
3. (No embedding) Within ‘no embedding’, as D increases, the proportion of problems that attained convergence_{opt}/convergence decreased for BARON and DIRECT. Surprisingly, for the KNITRO’s multi-start method, the proportion of solved problems is maintained, but the number of function evaluations and time increased dramatically.

Additional experiments. To see how robust REGO is to the changes in the parameters, we conduct three more experiments presented and discussed in Appendix C. In the first experiment, assuming that μ is bounded by \sqrt{D} (see page 50 for an explanation for this choice), we set δ to $\sqrt{D\bar{\delta}}$ for $\bar{\delta}$ chosen as in the main experiment. The second experiment tests REGO for four different values of d while keeping δ fixed and the third experiment tests REGO for three different values of δ keeping d fixed. In all three experiments, REGO performs well, particularly for BARON and KNITRO, solving most of the problems and exhibiting similar trends as in the main experiment.

6. Conclusions and future work

We study a general algorithmic framework for functions with low effective dimensionality that solves the reduced problem (RP) using a single Gaussian random embedding and a(ny) general global optimization solver. Our precise theoretical findings backed by the numerical experiments show that the success of (RP) is essentially independent of D and mainly depends on the gap between the embedding dimension d and the dimension of effective subspace d_e , and the ratio between the size of \mathcal{Y} (namely δ) and μ (the Euclidean distance to the closest affine subspace of minimizers). REGO with three standard global solvers produced high frequencies of convergence, generally outperforming the respective solver’s performance when applied directly to the problems (without the dimensionality reduction) in terms of proportion of problems solved and/or computational cost.

Our in-depth investigations are conceptual in nature, and there is clearly more work that needs to be done to make this framework practically applicable to global optimization problems with special structure. In particular, our REGO approach depends on knowing (an upper bound d on) the effective dimension d_e and the distance to the closest minimizer μ . As discussed on page 50 (in *Choosing d , μ and δ in practice* paragraph), future work may include estimating d_e prior to optimizing, noting that REGO does not need to learn the entire effective subspace only its dimension. Estimating d or d_e numerically is another possibility, as proposed in [43], where d is gradually increased until no significant changes in the best function value found are observed. Our theoretical choices for δ also depend on μ , which again needs estimating. In this case, choosing a box domain for f would provide a rough estimate for μ (as discussed on page 50), with the remark that REGO cannot (yet) guarantee feasibility with respect to given bounds. To achieve the latter, one needs to either add projection operators as in [4, 5, 49] or include the problem constraints in the formulation of (RP) and allow multiple random embeddings as in [38]. We investigated the latter approach for a constrained formulation of (P) over a box domain $\mathcal{X} = [-1, 1]^D$ in [7], where we introduce a new algorithmic framework X-REGO based on a reduced problem that is defined in terms of linear constraints $A\mathbf{y} \in \mathcal{X}$ instead of $\mathbf{y} \in [-\delta, \delta]^d$. This formulation ensures that the solutions of the reduced problem are always feasible with respect to the original bound-constrained problem. Furthermore, as the constraints of the X-REGO’s reduced problem do not involve δ , it eliminates the need to quantify this parameter (and thus estimate μ). In [7], we provide a thorough

theoretical analysis discovering relationships between the parameters of the problem (such as d and D) and the success of X-REGO. The analysis is based on the results derived in this paper.

Another potential direction is the use of other random matrix ensembles than Gaussian to generate the random embeddings. A first case to consider would be sub-Gaussian matrices, so that as many properties of the Gaussian case as possible, are inherited. For practical purposes, ensembles with finite support and/or sparse are of particular interest. For example, as mentioned in the Introduction, Naeybi *et al.* [34] use hashing matrices (one non-zero element per column) to generate random subspaces, which is computationally beneficial. However, we note that our theoretical results are heavily dependent on properties of Gaussian matrices and that deriving similar results for other matrix ensembles could be more challenging.

Lastly, real-life problems are often only approximately low dimensional and so their optimization requires further extensions and analysis of the random embedding framework. Nonetheless, we note that the techniques here, namely of random embeddings can still be applied in the approximate case, with some practical choice of d based on available computational budget; one can also increase d incrementally until no significant progress in the objective function value is obtained as suggested above. The latter was done, for example, in Li *et al.* [31], when training neural networks; they have observed that one low-dimensional random embedding seems sufficient to achieve, for instance, a high validation accuracy (see also [20, 38]).

Data availability statement

The data underlying this article are available in the article and in its online supplementary material.

Funding

The Alan Turing Institute (EP/N510129/1 and the Turing Project Scheme).

REFERENCES

1. BEN SALEM, M., BACHOC, F., ROUSTANT, O., GAMBOA, F. & TOMASO, L. (2019) Sequential dimension reduction for learning features of expensive black-box functions. <https://hal.archives-ouvertes.fr/hal-01688329/file/main.pdf>.
2. BERGSTRA, J. & BENGIO, Y. (2012) Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, **13**, 281–305.
3. BERNARDO, J. M. & SMITH, A. F. M. (2000) *Bayesian Theory*. Chichester: Wiley.
4. BINOIS, M., GINSBOURGER, D. & ROUSTANT, O. (2014) A warped kernel improving robustness in Bayesian optimization via random embeddings. Learning and Intelligent Optimization: 9th International Conference, LION 9. Revised Selected Papers, **8994**.
5. BINOIS, M., GINSBOURGER, D. & ROUSTANT, O. (2020) On the choice of the low-dimensional domain for global optimization via random embeddings. *Journal of Global Optimization*, **76**, 69–90.
6. BYRD, R. H., NOCEDAL, J. & WALTZ, R. A. (2006) Knitro: an integrated package for nonlinear optimization. Large-Scale Nonlinear Optimization. Boston, MA: Springer, pp. 35–59.
7. CARTIS, C., MASSART, E. & OTEMISSOV, A. (2020) Constrained global optimization of functions with low effective dimensionality using multiple random embeddings. arXiv e-prints, page arXiv:2009.10446.
8. CHEN, B., KRAUSE, A. & CASTRO, R. M. (2012) Joint optimization and variable selection of high-dimensional Gaussian processes. *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*. Madison, WI, USA: Omnipress, pp. 1423–1430.
9. CONSTANTINE, P. (2015) *Active Subspaces*. Philadelphia, PA: SIAM.

10. DAVIDSON, K. R. & SZAREK, S. (2001) Local operator theory, random matrices and Banach spaces. *Handbook on the Geometry of Banach Spaces*, Amsterdam, The Netherlands: Elsevier Science vol. 1, pp. 317–366.
11. DEMO, N., TEZZELE, M. & ROZZA, G. (2020) A supervised learning approach involving active subspaces for an efficient genetic algorithm in high-dimensional optimization problems. arXiv e-prints, arXiv:2006.07282.
12. DIXON, L. C. W. & SZEGÖ, G. P. (1975) *Towards Global Optimization*. New York: Elsevier.
13. DJOLONGA, J., KRAUSE, A. & CEVHER, V. (2013) High-dimensional Gaussian process bandits. *Proceedings of the 26th International Conference on Neural Information Processing Systems, NIPS'13*, Red Hook, NY, USA: Curran Associates Inc., pp. 1025–1033.
14. EDELMAN, A. (1988) Eigenvalues and condition numbers of random matrices. *SIAM J. Matrix Anal. Appl.*, **9**, 543–560.
15. ERIKSSON, D., DONG, K., LEE, E. H., BINDEL, D. & WILSON, A. G. (2018) Scaling Gaussian process regression with derivatives. *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, Red Hook, NY, USA: Curran Associates Inc., pp. 6868–6878.
16. ERNESTO, P. A. & DILIMAN, U. P. (2005) MVF—multivariate test functions library in C for unconstrained global. *Optimization*. <http://www.geocities.ws/eadorio/mvf.pdf>
17. FANG, K., KOTZ, S. & NG, K. W. (1990) *Symmetric Multivariate and Related Distributions*. London: Chapman and Hall.
18. FINKEL, D. (2003) Direct optimization algorithm user guide. *Technical Report CRSC-TR03-11*. North Carolina State University, Center for Research in Scientific Computation. Raleigh, North Carolina. <http://www.ncsu.edu/crsc/reports/ftp/pdf/crsc-tr03-11.pdf>.
19. FORNASIER, M., SCHNASS, K. & VYBIRAL, J. (2012) Learning functions of few arbitrary linear parameters in high dimensions. *Found. Comput. Math.*, **12**, 229–262.
20. FRANKLE, J. & CARBIN, M. (2018) The lottery ticket hypothesis: finding sparse, trainable neural networks. International Conference on Learning Representations.
21. GABLONSKY, J. M. & KELLEY, C. T. (2001) A locally-biased form of the direct algorithm. *J. Global Optim.*, **21**, 27–37.
22. GARNETT, R., OSBORNE, M. A. & HENNIG, P. (2014) Active learning of linear embeddings for Gaussian processes. *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, UAI'14*, Arlington, Virginia, USA: AUAI Press, pp. 230–239.
23. GAVANA, A. (2018) Global optimization benchmarks and AMPGO. http://infinity77.net/global_optimization/.
24. GUPTA, A. K. & NAGAR, D. K. (2000) *Matrix Variate Distributions*. New York: Chapman and Hall/CRC.
25. GUPTA, A. K. & SONG, D. (1997) Lp-Norm spherical distribution. *J. Statist. Plann. Inference*, **60**, 241–260.
26. HUTTER, F., HOOS, H. & LEYTON-BROWN, K. (2014) An efficient approach for assessing hyperparameter importance. *Proceedings of the 31st International Conference on International Conference on Machine Learning, ICML'14*, vol. 32, Beijing, China: PMLR, pp. 754–762.
27. JONES, D. R., PERTTUNEN, C. D. & STUCKMAN, B. E. (1993) Lipschitzian optimization without the Lipschitz constant. *J. Optim. Theory Appl.*, **79**, 157–181.
28. KANDASAMY, K., SCHNEIDER, J. & PÓCZOS, B. (2015) High dimensional Bayesian optimisation and bandits via additive models. *Proceedings of the 32nd International Conference on International Conference on Machine Learning, ICML'15*, vol. 37, Lille, France: PMLR, pp. 295–304.
29. KNIGHT, C. G., KNIGHT, S. H. E., MASSEY, N., AINA, T., CHRISTENSEN, C., FRAME, D. J., KETTLEBOROUGH, J. A., MARTIN, A., PASCOE, S., SANDERSON, B., STAINFORTH, D. A. & ALLEN, M. R. (2007) Association of parameter, software, and hardware variation with large-scale behavior across 57,000 climate models. *Proc. Natl. Acad. Sci. USA*, **104**, 12259–12264.
30. LEE, P. M. (2012) *Bayesian Statistics: An Introduction*, 4th edn. Chichester: Wiley.
31. LI, C., FARKHOOR, H., LIU, R. & YOSINSKI, J. (2018) Measuring the intrinsic dimension of objective landscapes. arXiv e-prints, arXiv:1804.08838.
32. LI, C.-L., KANDASAMY, K., PÓCZOS, B. & SCHNEIDER, J. (2016) High dimensional Bayesian optimization via restricted projection pursuit models. *Proceedings of the 19th International Conference on Artificial*

- Intelligence and Statistics*, vol. 51. Proceedings of Machine Learning Research, Cadiz, Spain: PMLR, pp. 884–892.
33. LUKACZYK, T. W., CONSTANTINE, P., PALACIOS, F. & ALONSO, J. J. (2014) Active subspaces for shape optimization. *10th AIAA Multidisciplinary Design Optimization Conference*. Reston, Virginia: American Institute of Aeronautics and Astronautics, Inc.
 34. NAYEBI, A., MUNTEANU, A. & POLOCZEK, M. (2019) A framework for Bayesian optimization in embedded subspaces. *Proceedings of the 36th International Conference on Machine Learning*, vol. 97. Proceedings of Machine Learning Research, Long Beach, California: PMLR, pp. 4752–4761.
 35. NEUMAIER, A., SHCHERBINA, O., HUYER, W. & VINKÓ, T. (2005) A comparison of complete global optimization solvers. *Math. Program.*, **103**, 335–356.
 36. NEUMAN, E. (2013) Inequalities and bounds for the incomplete gamma function. *Results Math.*, **63**, 1209–1214.
 37. NIST/SEMATECH (2018) E-handbook of statistical methods. <https://www.itl.nist.gov/div898/handbook/eda/section3/eda3666.htm>.
 38. QIAN, H., HU, Y.-Q. & YU, Y. (2016) Derivative-free optimization of high-dimensional non-convex functions by sequential random embeddings. *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI'16*, Palo Alto, California: AAAI Press, pp. 1946–1952.
 39. QIAN, H. & YU, Y. (2017) Solving high-dimensional multi-objective optimization problems with low effective dimensions. *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI'17*, San Francisco, California: AAAI Press, pp. 875–881.
 40. ROLLAND, P., SCARLETT, J., BOGUNOVIC, I. & CEVHER, V. (2018) High-dimensional Bayesian optimization via additive models with overlapping groups. *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, **84**, pp. 298–307.
 41. RUDELSON, M. & VERSHYNIN, R. (2009) Smallest singular value of a random rectangular matrix. *Comm. Pure Appl. Math.*, **62**, 1707–1739.
 42. SAHINIDIS, N. V. (2014) BARON 14.3.1: global optimization of mixed-integer nonlinear programs, user's manual.
 43. SANYANG, M. L. & KABÁN, A. (2016) REMEDA: random embedding EDA for optimising functions with intrinsic dimension. *International Conference on Parallel Problem Solving from Nature, PPSN XIV*, Palo Cham: Springer International Publishing, pp. 859–868.
 44. SURJANOVIC, S. & BINGHAM, D. (2013) Virtual library of simulation experiments: test functions and datasets. <https://www.sfu.ca/ssurjano/>.
 45. TAWARMALANI, M. & SAHINIDIS, N. V. (2005) A polyhedral branch-and-cut approach to global optimization. *Math. Program.*, **103**, 225–249.
 46. TYAGI, H. & CEVHER, V. (2014) Learning non-parametric basis independent models from point queries via low-rank methods. *Appl. Comput. Harmon. Anal.*, **37**, 389–412.
 47. VERSHYNIN, R. (2018) *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.
 48. WANG, Z., GEHRING, C., KOHLI, P. & JEGELKA, S. (2018) Batched large-scale Bayesian optimization in high-dimensional spaces. *International Conference on Artificial Intelligence and Statistics (AISTATS)*. New York, United States: Association for Computing Machinery.
 49. WANG, Z., HUTTER, F., ZOGHI, M., MATHESON, D. & DE FREITAS, N. (2016) Bayesian optimization in a billion dimensions via random embeddings. *J. Artificial Intelligence Res.*, **55**, 361–387.
 50. WHEEDEN, R. L. (2015) *Measure and Integral: An Introduction to Real Analysis*, 2nd edn. Boca Raton: Chapman and Hall/CRC.

A. Technical definitions and results

A.1 Gaussian random matrices

DEFINITION A.1 (Gaussian matrix). A Gaussian (random) matrix is a matrix whose each entry is an independent standard normal random variable.

Gaussian matrices have been well studied with many results available at hand. Here, we mention a few key properties of Gaussian matrices that we use in the analysis; for a collection of results pertaining to Gaussian matrices and other related distributions, refer to [24, 47].

Gaussian random matrices are known to be invariant with respect to orthogonal transformations.

THEOREM A.2 (see [24, Theorem 2.3.10]) Let \mathbf{A} be an $D \times d$ Gaussian random matrix. If $\mathbf{U} \in \mathbb{R}^{D \times p}$, $D \geq p$ and $\mathbf{V} \in \mathbb{R}^{d \times q}$, $d \geq q$ are orthogonal, then $\mathbf{U}^T \mathbf{A} \mathbf{V}$ is a Gaussian random matrix.

A related notion that plays an important role in the study of Gaussian matrices is the Wishart distribution represented by matrix $\mathbf{A}^T \mathbf{A}$ (or $\mathbf{A} \mathbf{A}^T$), where \mathbf{A} is an overdetermined (underdetermined) Gaussian matrix. A Wishart matrix is positive definite with probability 1.

THEOREM A.3 (see [24, Theorem 3.2.1]) Let \mathbf{A} be an $D \times d$ Gaussian random matrix, $D \geq d$. Then, the Wishart matrix $\mathbf{A}^T \mathbf{A}$ is positive definite with probability 1.

The immediate consequence of the result is that the Wishart matrix is non-singular with probability 1.

A.2 Chi-squared random variable

DEFINITION A.4 (Chi-squared random variable). Given a collection Z_1, Z_2, \dots, Z_n of n independent standard normal variables, a random variable $X = Z_1^2 + Z_2^2 + \dots + Z_n^2$ is said to follow the chi-squared distribution with n degrees of freedom. We denote this by $X \sim \chi_n^2$.

The following lemma provides a notable relationship between the inverse of the Wishart matrix and the chi-squared random variable.

LEMMA A.5 (see [24, Corollary 3.3.13.1.]) Let \mathbf{A} be an $D \times d$ Gaussian matrix, $D \geq d$, and let $\mathbf{z} \in \mathbb{R}^d$ be a fixed non-zero vector. Then,

$$\frac{\|\mathbf{z}\|^2}{\mathbf{z}^T (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{z}} \sim \chi_{D-d+1}^2.$$

In the following lemma, we derive an upper bound for the cumulative density function (c.d.f.) of the chi-squared random variable.

LEMMA A.6 Let $X \sim \chi_n^2$. Then,

$$\mathbb{P}[X \leq x] \leq \frac{4}{n(n+2)\Gamma(n/2)} \left(1 + \frac{n}{2} e^{-x/2}\right) (x/2)^{n/2}.$$

Proof. Recall the c.d.f. of the chi-square random variable (see, e.g. [37]):

$$\mathbb{P}[X \leq x] = \frac{\gamma(n/2, x/2)}{\Gamma(n/2)}$$

for $x > 0$, where $\gamma(n/2, x/2)$ is the lower incomplete gamma function (see [36]) defined as

$$\gamma(n/2, x/2) = \int_0^{x/2} u^{n/2-1} e^{-u} du.$$

We obtain the desired result by applying the following upper bound on $\gamma(n/2, x/2)$ (see [36, Theorem 4.1]):

$$\gamma(n/2, x/2) \leq \frac{4}{n(n+2)} \left(1 + \frac{n}{2} e^{-x/2} \right) (x/2)^{n/2}.$$

□

A.3 The inverse chi-squared random variable

DEFINITION A.7 (Inverse chi-squared random variable). Given $X \sim \chi_n^2$, a random variable $Y = 1/X$ is said to follow the inverse chi-squared distribution with n degrees of freedom. We denote this by $Y \sim 1/\chi_n^2$ [30, A5].

LEMMA A.8 (see [30, A5]) Let $Y \sim 1/\chi_n^2$ and $W = sY$ for a positive real s . Then,

$$\mathbb{E}[W] = \frac{s}{n-2}$$

provided that $n > 2$.

LEMMA A.9 Let Y and R be two random variables such that $Y \sim 1/\chi_n^2$ and $R = \sqrt{Y}$. Then, the probability density function (p.d.f.) $g(\hat{r})$ of R is given by

$$g(\hat{r}) = \frac{2^{-n/2+1}}{\Gamma(n/2)} \hat{r}^{-n-1} e^{-1/(2\hat{r}^2)}.$$

Proof. The p.d.f. $g(\hat{r})$ of R satisfies

$$g(\hat{r}) = \frac{d}{d\hat{r}} \mathbb{P}[\sqrt{Y} < \hat{r}] = \frac{d}{d\hat{r}} \mathbb{P}[Y < \hat{r}^2] = 2\hat{r}h(n, \hat{r}^2), \quad (\text{A1})$$

where $h(n, \cdot)$ denotes the p.d.f. of the inverse chi-squared random variable Y given by (see, e.g. [30, A5])

$$h(n, y) = \frac{1}{2^{n/2} \Gamma(n/2)} y^{-n/2-1} e^{-1/(2y)} \text{ for } y > 0.$$

□

A.4 Spherically distributed random vectors

DEFINITION A.10 An $D \times 1$ random vector \mathbf{x} is said to have a spherical distribution if for every orthogonal $D \times D$ matrix \mathbf{U} ,

$$\mathbf{U}\mathbf{x} \stackrel{\text{law}}{=} \mathbf{x}.$$

Below are some useful facts about symmetrically distributed random vectors.

LEMMA A.11 [17, p. 13] Let \mathbf{x} and \mathbf{y} be random vectors such that $\mathbf{x} \stackrel{\text{law}}{=} \mathbf{y}$, and let $f_i(\cdot)$, $i = 1, 2, \dots, m$, be measurable functions. Then,

$$(f_1(\mathbf{x}) \ f_2(\mathbf{x}) \ \dots \ f_m(\mathbf{x}))^T \stackrel{\text{law}}{=} (f_1(\mathbf{y}) \ f_2(\mathbf{y}) \ \dots \ f_m(\mathbf{y}))^T.$$

LEMMA A.12 ([17, Corollary, p. 30]) If $D \times 1$ random vector \mathbf{x} has a spherical distribution, then

$$\mathbf{x} \stackrel{\text{law}}{=} r\mathbf{u},$$

where \mathbf{u} is distributed uniformly on the $(D-1)$ -dimensional unit sphere and r is a univariate random variable independent of \mathbf{u} .

THEOREM A.13 ([17, Theorem 2.3]) Let $\mathbf{x} \stackrel{\text{law}}{=} \mathbf{r}\mathbf{u}$ be a spherically distributed $D \times 1$ random vector with $\mathbb{P}[\mathbf{x} = \mathbf{0}] = 0$. Then,

$$\|\mathbf{x}\| \stackrel{\text{law}}{=} r \text{ and } \|\mathbf{x}\|^{-1}\mathbf{x} \stackrel{\text{law}}{=} \mathbf{u}.$$

Moreover, $\|\mathbf{x}\|$ and $\|\mathbf{x}\|^{-1}\mathbf{x}$ are independent.

THEOREM A.14 (see [25, Theorem 2.1.]) Let $\mathbf{x} \stackrel{\text{law}}{=} \mathbf{r}\mathbf{u}$ be a spherically distributed $D \times 1$ random vector with $\mathbb{P}[\mathbf{x} = \mathbf{0}] = 0$, where r is independent of \mathbf{u} with p.d.f. $h(\cdot)$. Then, p.d.f. $g(\hat{\mathbf{x}})$ of \mathbf{x} is given by

$$g(\hat{\mathbf{x}}) = \frac{\Gamma(D/2)}{2\pi^{D/2}} h(\|\hat{\mathbf{x}}\|) \|\hat{\mathbf{x}}\|^{1-D}.$$

For more details regarding spherical distributions refer to [3, 17, 24].

A.5 The least Euclidean norm solution to the random linear system

The present section establishes key properties of the least Euclidean norm solution to the underdetermined random linear system.

- (C) Let $\bar{\mathbf{B}}$ be a $d_e \times d$ Gaussian matrix, where $d_e \leq d$, and let $\mathbf{z} \in \mathbb{R}^{d_e}$ be a fixed non-zero vector. Denote by \mathbf{y}_2 the least 2-norm solution to $\bar{\mathbf{B}}\mathbf{y} = \mathbf{z}$.

LEMMA A.15 Given (C), \mathbf{y}_2 satisfies

$$\frac{\|\mathbf{z}\|_2^2}{\|\mathbf{y}_2\|_2^2} \sim \chi_{d-d_e+1}^2.$$

Proof. The least Euclidean norm solution \mathbf{y}_2 to $\bar{\mathbf{B}}\mathbf{y} = \mathbf{z}$ is given by

$$\mathbf{y}_2 = \bar{\mathbf{B}}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z}.$$

For its Euclidean norm, we have

$$\begin{aligned} \|\mathbf{y}_2\|_2^2 &= (\bar{\mathbf{B}}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z})^T \bar{\mathbf{B}}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z} \\ &= \mathbf{z}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z}. \end{aligned}$$

Using Lemma A.5, we obtain the desired result:

$$\frac{\|\mathbf{z}\|_2^2}{\|\mathbf{y}_2\|_2^2} = \frac{\|\mathbf{z}\|^2}{\mathbf{z}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z}} \sim \chi_{d-d_e+1}^2.$$

□

LEMMA A.16 Given (C), \mathbf{y}_2 follows a spherical distribution.

Proof. Let \mathbf{S} be any $d \times d$ orthogonal matrix. Let $f : \mathbb{R}^{d_e d \times 1} \rightarrow \mathbb{R}^{d \times 1}$ be a vector-valued function defined as

$$f(\text{vec}(\bar{\mathbf{B}})) = \bar{\mathbf{B}}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z},$$

where $\text{vec}(\bar{\mathbf{B}})$ denotes the $Dd \times 1$ vector $(\bar{\mathbf{b}}_1^T \bar{\mathbf{b}}_2^T \cdots \bar{\mathbf{b}}_d^T)^T$ with $\bar{\mathbf{b}}_i$ being the i th column vector of $\bar{\mathbf{B}}$. Using the fact that the inverse of a matrix is equal to the ratio of its adjugate to its determinant we can express

f as

$$f(\text{vec}(\bar{\mathbf{B}})) = \left(\frac{p_1(\bar{\mathbf{B}})}{q(\bar{\mathbf{B}})} \quad \frac{p_2(\bar{\mathbf{B}})}{q(\bar{\mathbf{B}})} \quad \dots \quad \frac{p_d(\bar{\mathbf{B}})}{q(\bar{\mathbf{B}})} \right)^T,$$

where $p_i(\bar{\mathbf{B}})$ for $1 \leq i \leq d$ are some polynomials of the entries of $\bar{\mathbf{B}}$ and $q(\bar{\mathbf{B}})$ is the determinant of $\bar{\mathbf{B}}\bar{\mathbf{B}}^T$.

We first would like to prove that f is a measurable function. Recall that a function is measurable if and only if each of its components is measurable. It is enough to show that p_1/q is measurable; the same argument will apply to the rest of its components. First, we note that

- (i) p_1 and q are measurable;
- (ii) q is non-zero almost everywhere.

To prove (i), observe that the polynomials p_1 and q are sums of scalar multiples of products of standard normal random variables, which by definition are measurable. Sums, scalar multiples and products of measurable functions are measurable; hence, p_1 and q must be measurable. To prove (ii), we refer to Theorem A.3, which says that the matrix $\bar{\mathbf{B}}\bar{\mathbf{B}}^T$ is positive definite with probability 1 implying that all of its eigenvalues are strictly positive with probability 1. Then, (ii) follows from the fact that the determinant of the symmetric square matrix is equal to the product of its eigenvalues. Now, we can apply [50, Theorem 4.10] to deduce that p_1/q is measurable; this completes the proof that f is measurable.

For $\mathbf{y}_2 = \bar{\mathbf{B}}^T (\bar{\mathbf{B}}\bar{\mathbf{B}}^T)^{-1} \mathbf{z}$, we have

$$\mathbf{y}_2 = f(\text{vec}(\bar{\mathbf{B}})) \text{ and } \mathbf{S}\mathbf{y}_2 = f(\text{vec}(\bar{\mathbf{B}}\mathbf{S}^T)).$$

According to Theorem A.2, $\text{vec}(\bar{\mathbf{B}}) \stackrel{\text{law}}{=} \text{vec}(\bar{\mathbf{B}}\mathbf{S}^T)$. Then, by applying Lemma A.11, we obtain

$$\mathbf{y}_2 \stackrel{\text{law}}{=} \mathbf{S}\mathbf{y}_2.$$

Hence, \mathbf{y}_2 follows a spherical distribution by Definition A.10. □

LEMMA A.17 Given (C), the probability density function of \mathbf{y}_2 is given by

$$g(\hat{\mathbf{y}}) = \pi^{-d/2} \left[\frac{\Gamma(d/2)}{\Gamma(n/2)} \right] \left(\frac{\|\mathbf{z}\|}{\sqrt{2}} \right)^n (\hat{\mathbf{y}}^T \hat{\mathbf{y}})^{-(n+d)/2} e^{-\|\mathbf{z}\|^2 / (2\hat{\mathbf{y}}^T \hat{\mathbf{y}})},$$

where $n = d - d_e + 1$.

Proof. The fact that \mathbf{y}_2 has a spherical distribution is a key ingredient in the proof. To simplify the derivations, let us assume for now that $\|\mathbf{z}\| = 1$.

By Lemma A.12,

$$\mathbf{y}_2 \stackrel{\text{law}}{=} r\mathbf{u},$$

where r is a univariate random variable and where \mathbf{u} is a random vector distributed uniformly on the $(d - 1)$ -dimensional unit sphere; moreover, r and \mathbf{u} are independent.

Our first goal is to show that r and $\|\mathbf{y}_2\|$ have the same distribution. This fact follows immediately from Theorem A.13 if we show that $\mathbb{P}[\mathbf{y}_2 = \mathbf{0}] = 0$. Let $W \sim 1/\chi_{d-d_e+1}^2$. According to Lemma A.15, we have

$$\|\mathbf{y}_2\|^2 \stackrel{\text{law}}{=} W, \tag{A2}$$

which we use in the second equation below:

$$\mathbb{P}[\mathbf{y}_2 = \mathbf{0}] = \mathbb{P}[\|\mathbf{y}_2\|^2 = 0] = \mathbb{P}[W = 0].$$

Since W is a continuous random variable, $\mathbb{P}[W = 0] = 0$. This proves that

$$r \stackrel{\text{law}}{=} \|y_2\|. \quad (\text{A3})$$

Combining (A2) and (A3), we conclude that r has the same distribution as $W^{1/2}$. The probability density function of $W^{1/2}$ —and, consequently, of r —derived in Lemma A.9 is given by

$$h(\hat{r}) = \frac{2^{-n/2+1}}{\Gamma(n/2)} \hat{r}^{-n-1} e^{-1/(2\hat{r}^2)}. \quad (\text{A4})$$

Theorem A.14 allows us to express the p.d.f. of y_2 in terms of the p.d.f. of r :

$$g(\hat{y}) = \frac{\Gamma(d/2)}{2\pi^{d/2}} (\hat{y}^T \hat{y})^{(1-d)/2} h(\sqrt{\hat{y}^T \hat{y}}).$$

By using (A4) for $h(\cdot)$ in the above, we obtain

$$g(\hat{y}) = 2^{-n/2} \pi^{-d/2} \frac{\Gamma(d/2)}{\Gamma(n/2)} (\hat{y}^T \hat{y})^{-(n+d)/2} e^{-1/(2\hat{y}^T \hat{y})}. \quad (\text{A5})$$

To derive the p.d.f. for arbitrary non-zero z , we consider the linear transformation $\bar{y} = \|z\|\hat{y}$. The Jacobian of the transformation is equal to $1/\|z\|^d$. Thus, the p.d.f. $\bar{g}(\bar{y})$ of \bar{y} satisfies

$$\bar{g}(\bar{y}) = \frac{g(\bar{y}/\|z\|)}{\|z\|^d},$$

which together with (A5) yields the desired result. \square

B. Problem set

Table B3 contains the explicit formula, domain and global minimum of the functions used to generate the high-dimensional test set. The problem set contains 19 problems taken from [16, 23, 44]. Problems that cannot be solved by BARON are marked with “*”. Problems that will not be solved by KNITRO are marked with “°”.

C. Additional experiments

We conducted three more experiments to test REGO’s robustness to changes in the parameters. In all three experiments, the same budget and termination criteria as in the main experiment are used.

- (A) In this experiment, we assume that no good estimate for μ is known and that μ can be as large as \sqrt{D} (for example, when $\mathcal{X} = [-1, 1]^D$ constraint is imposed). We test REGO with the following parameters: $(d_e, 8.0 \times \sqrt{D})$, $(d_e + 1, 2.2 \times \sqrt{D})$, $(d_e + 2, 1.3 \times \sqrt{D})$ and $(d_e + 3, 1.0 \times \sqrt{D})$. Results are presented in Fig. C8 in Appendix C.
- (B) We fix δ to be $7.5\sqrt{d_e}$ and vary d . The following parameters are used: $(d_e, 7.5 \times \sqrt{d_e})$, $(d_e + 1, 7.5 \times \sqrt{d_e})$, $(d_e + 2, 7.5 \times \sqrt{d_e})$ and $(d_e + 3, 7.5 \times \sqrt{d_e})$. Results are presented in Fig. C9 in Appendix C.
- (C) We fix $d = d_e + 1$ and vary δ ($= 5\sqrt{d_e}, 7.5\sqrt{d_e}, 10\sqrt{d_e}$). The following parameters are used: $(d_e + 1, 5 \times \sqrt{d_e})$, $(d_e + 1, 7.5 \times \sqrt{d_e})$ and $(d_e + 1, 10 \times \sqrt{d_e})$. In the figures, we also include curves for $\delta_{opt} = 2.2\sqrt{d_e}$ taken from the main experiment. Results are presented in Fig. C10 in Appendix C.

TABLE B3 *The problem set listed in alphabetical order*

| Function | Domain | Global minima |
|--------------------------|--------------------------------------|-------------------------------|
| 1) Beale [16] | $\mathbf{x} \in [-4.5, 4.5]^2$ | $g(\mathbf{x}^*) = 0$ |
| 2) *Branin [16] | $x_1 \in [-5, 10], x_2 \in [0, 15]$ | $g(\mathbf{x}^*) = 0.397887$ |
| 3) Brent [23] | $\mathbf{x} \in [-10, 10]^2$ | $g(\mathbf{x}^*) = 0$ |
| 4) °Bukin N.6 [44] | $x_1 \in [-15, -5], x_2 \in [-3, 3]$ | $g(\mathbf{x}^*) = 0$ |
| 5) *Easom [16] | $\mathbf{x} \in [-100, 100]^2$ | $g(\mathbf{x}^*) = -1$ |
| 6) Goldstein-Price [16] | $\mathbf{x} \in [-2, 2]^2$ | $g(\mathbf{x}^*) = 3$ |
| 7) Hartmann 3 [16] | $\mathbf{x} \in [0, 1]^3$ | $g(\mathbf{x}^*) = -3.86278$ |
| 8) Hartmann 6 [16] | $\mathbf{x} \in [0, 1]^6$ | $g(\mathbf{x}^*) = -3.32237$ |
| 9) *Levy [44] | $\mathbf{x} \in [-10, 10]^4$ | $g(\mathbf{x}^*) = 0$ |
| 10) Perm 4, 0.5 [44] | $\mathbf{x} \in [-4, 4]^4$ | $g(\mathbf{x}^*) = 0$ |
| 11) Rosenbrock [44] | $\mathbf{x} \in [-5, 10]^3$ | $g(\mathbf{x}^*) = 0$ |
| 12) Shekel 5 [44] | $\mathbf{x} \in [0, 10]^4$ | $g(\mathbf{x}^*) = -10.1532$ |
| 13) Shekel 7 [44] | $\mathbf{x} \in [0, 10]^4$ | $g(\mathbf{x}^*) = -10.4029$ |
| 14) Shekel 10 [44] | $\mathbf{x} \in [0, 10]^4$ | $g(\mathbf{x}^*) = -10.5364$ |
| 15) *Shubert [44] | $\mathbf{x} \in [-10, 10]^2$ | $g(\mathbf{x}^*) = -186.7309$ |
| 16) Six-hump camel [44] | $x_1 \in [-3, 3], x_2 \in [-2, 2]$ | $g(\mathbf{x}^*) = -1.0316$ |
| 17) Styblinski-Tang [44] | $\mathbf{x} \in [-5, 5]^4$ | $g(\mathbf{x}^*) = -156.664$ |
| 18) Trid [44] | $\mathbf{x} \in [-25, 25]^5$ | $g(\mathbf{x}^*) = -30$ |
| 19) Zettl [16] | $\mathbf{x} \in [-5, 5]^2$ | $g(\mathbf{x}^*) = -0.00379$ |

Conclusions

- (A) We test robustness of REGO assuming that μ is equal to \sqrt{D} (which makes δ to be relatively large and dependent on D). Despite this dependence, the frequency of convergence for BARON and KNITRO is high showing mild dependence on D .
- (B) The purpose of this experiment is to see how different values of d affect the performance of REGO while δ is kept constant. For larger d , we expect (RP) to be successful with higher chance. Nonetheless, the results show that sometimes, for larger d , REGO's performance may be compromised; this is for example true for BARON's convergence_{opt}. Since δ is set to a relatively large value, (RP) is successful with high probability even for smallest d . This suggests that as long as d and δ produce relatively high chance of success of (RP), one should stop increasing their values lest convergence to the global minimum require larger computational resources.
- (C) In this experiment, we apply REGO with different values of δ while keeping d constant. The results display no significant differences between the performances with different parameters. Even the results with the optimal δ (used in the main experiment) do not differ considerably from the one with the largest δ except for BARON where the former wins in terms of convergence_{opt} and CPU time. The results of this experiment together with the results in (B) indicate that it is better to increase δ and keep d constant if one wants to increase success of (RP) with minimal increase in computational cost.

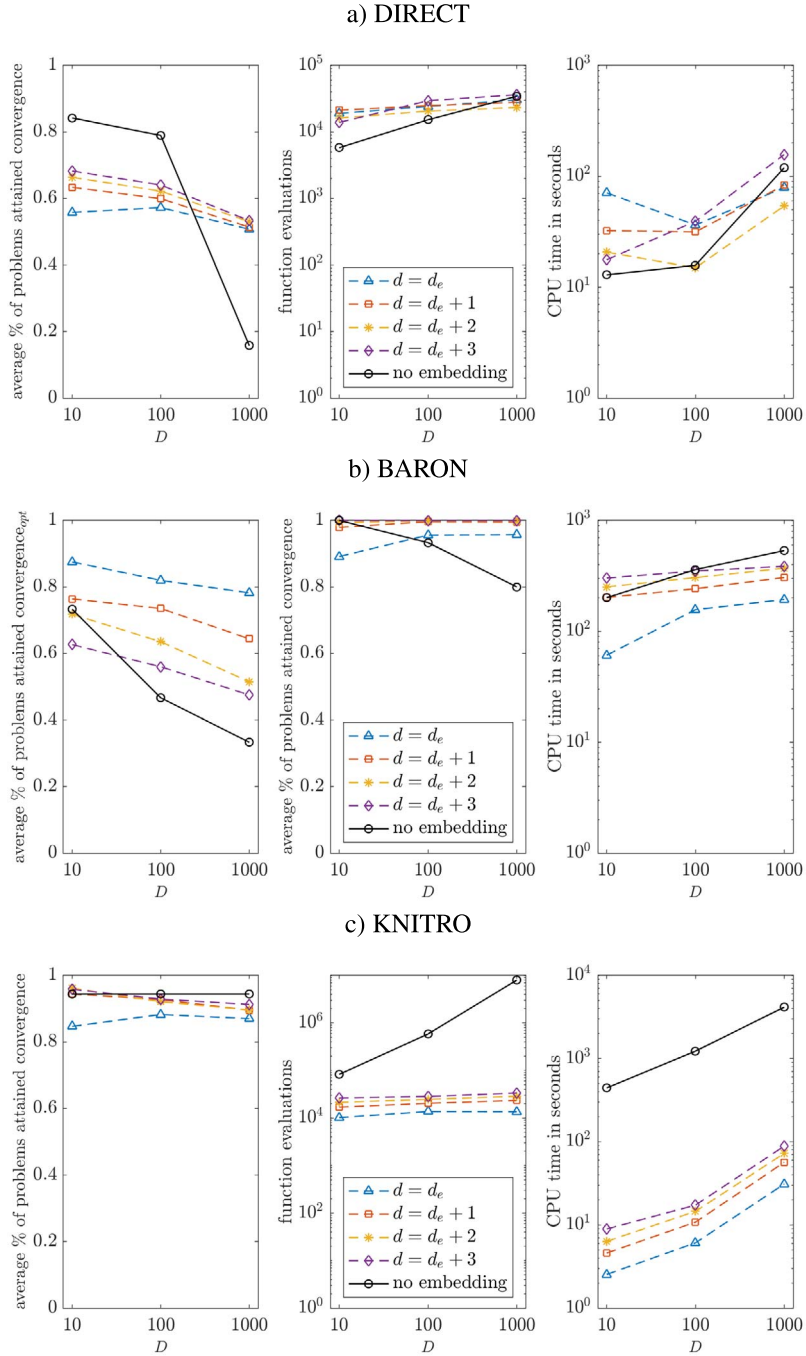


FIG. C8. Experiment A: REGO with DIRECT, BARON and KNITRO with $(d, \delta) = (d_e, 8.0 \times \sqrt{D})$, $(d_e + 1, 2.2 \times \sqrt{D})$, $(d_e + 2, 1.3 \times \sqrt{D})$ and $(d_e + 3, 1.0 \times \sqrt{D})$.

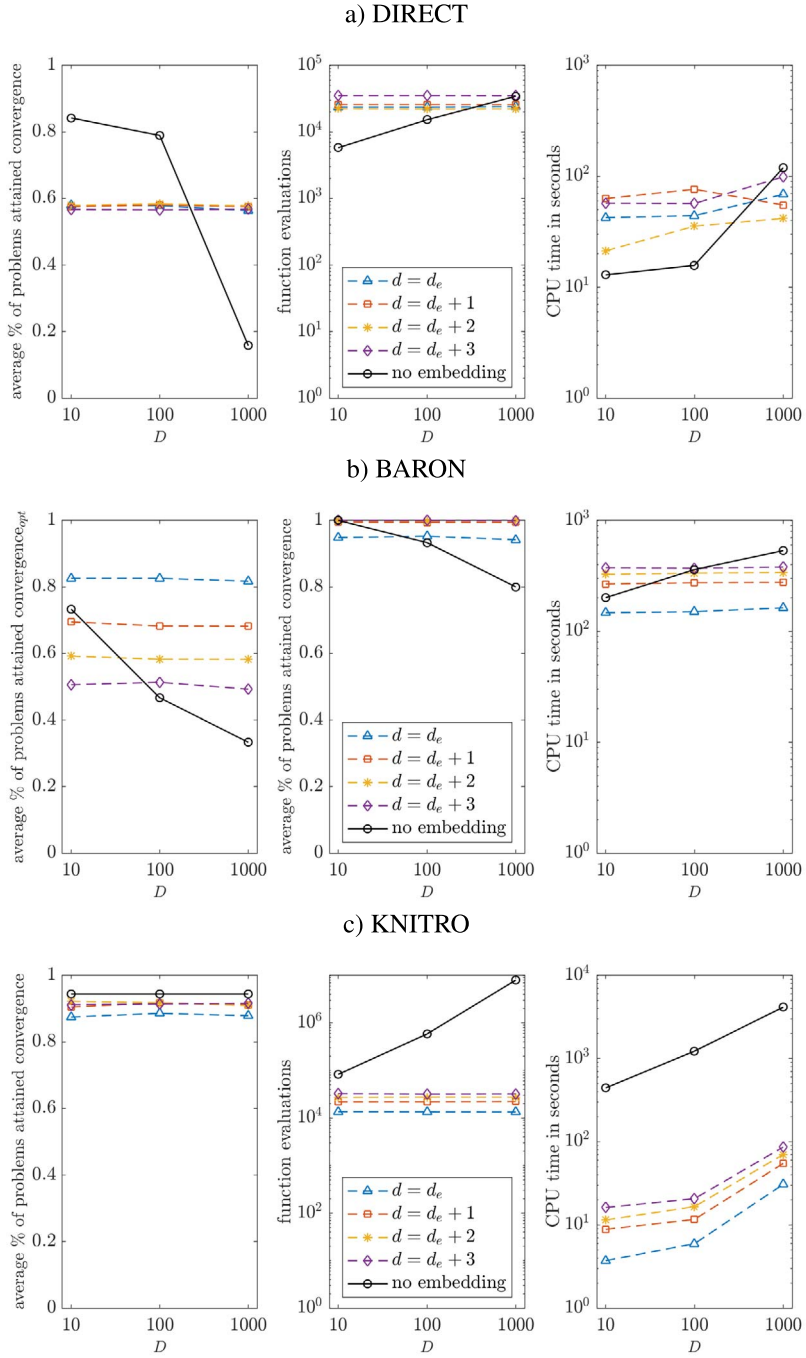


FIG. C9. Experiment B: REGO with DIRECT, BARON and KNITRO with $\delta = 7.5\sqrt{d_e}$ fixed and $d = d_e, d_e + 1, d_e + 2$ and $d_e + 3$.

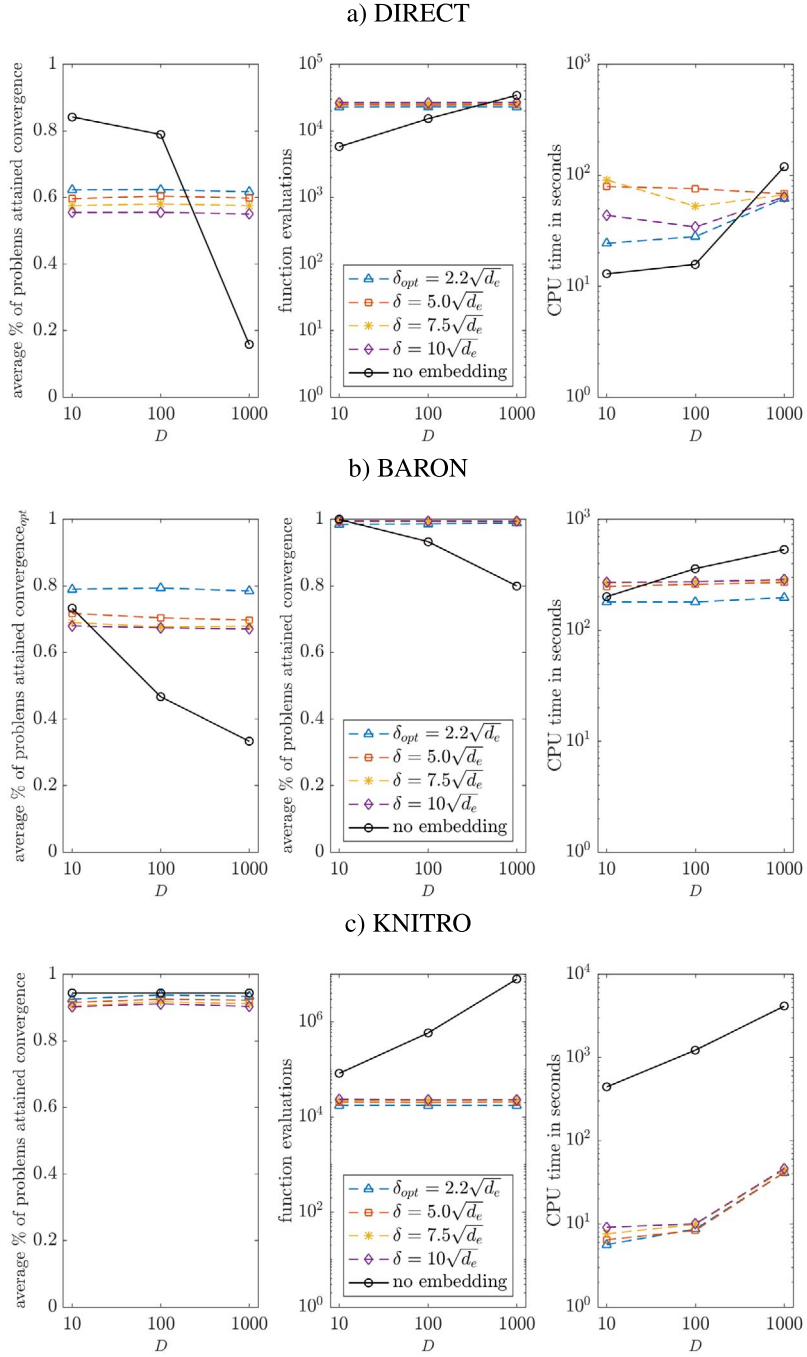


FIG. C10. Experiment C: REGO with DIRECT, BARON and KNITRO with $d = d_e + 1$ fixed and $\delta = 5\sqrt{d_e}, 7.5\sqrt{d_e}, 10\sqrt{d_e}$ and $2.2\sqrt{d_e}(\delta_{opt})$.