

On a Population Sizing Model for Evolution Strategies in Multimodal Landscapes

Lisa Schönerberger¹ and Hans-Georg Beyer²

Abstract—This article derives a population sizing model for standard evolution strategies (ES) in highly multimodal fitness landscapes with exponentially many local optima. The Rastrigin, Bohachevsky, and Ackley test functions are considered. Due to the highly nonconvex structure of these functions a detailed analytical description of the behavior of the ES is a challenge. Therefore, a model is derived that simplifies the complex structure of the functions under consideration. The main idea of this model is the interpretation of local landscape oscillations as frozen noise. This allows for an estimation of the success probability of the ES converging to the global optimum and in turn an estimation of the population size required. It is shown that the population size scales usually sublinearly with the search space dimension N . For the Rastrigin and Bohachevsky function, the population size scales with $\mathcal{O}(\sqrt{N} \ln(N))$. As for Ackley, the scaling behavior depends strongly on the initial values. If the algorithm starts in a certain vicinity of the global optimizer, the dependence on the dimension N is rather weak. However, if the initial value exceeds a certain distance R to the optimizer, the population size scales exponentially with R .

Index Terms—Evolution strategies (ES), global convergence, global optimization, multimodal objective function, population sizing.

I. INTRODUCTION

OPTIMIZATION in highly multimodal, real-valued fitness landscapes is a challenging topic and the theoretical analysis is still in its infancy. Due to their underlying stochastic nature, evolution strategies (ES) have proven to be well suited for the optimization of highly multimodal problems. However, the success rate of ES for finding the global optimizer depends strongly on the population size, as already observed in [1] for the CMA-ES [2]. This also holds for $(\mu/\mu_I, \lambda)$ -ES using isotropic mutations with σ self-adaptation (σ SA) or cumulative step-size adaptation (CSA) for mutation strength control. On the one hand the population size must not be too small; in this case the ES will not find the global optimizer. A too large population size on the other hand would require too many resources. Furthermore, the optimal population size depends

on the search space dimensionality N . This is because the number of local minima typically increases with increasing N .

This increase usually disqualifies the application of classical nonlinear optimization algorithms. That is, single-run gradient-based strategies are not able to find the global optimizer. Instead, restart strategies must be applied. However, since the number of restarts must increase exponentially with N , such strategies are not a real alternative because this implies an exponentially increasing number of function evaluations. In contrast, there are indications that for ES the population size scales slower than exponentially, at most quadratically or even partially sublinearly with the problem dimension, thus, keeping the number of function evaluations at a moderate order. In [1], this was experimentally demonstrated for several multimodal test functions, such as Rastrigin, Bohachevsky, Ackley, Schaffer, and Schwefel. In [3], these experimental results were confirmed by an analytical model for the Rastrigin function.

It is the aim of this article to extend this model to other multimodal test functions that share certain similarities with the Rastrigin function, i.e., the Bohachevsky and the Ackley function. As a result, population sizing equations will be obtained that scale sublinearly with the search space dimension N . The approach taken and findings made may be regarded as one of the first steps toward an analysis of ES on highly multimodal optimization problems. The population sizing results can be used to evaluate the effectiveness of population sizing rules. This in turn may be a starting point for further algorithm as well as benchmark designs. Last but not least, the analysis also provides a deeper understanding of how ES finds global optima.

The remainder of this article is organized as follows. After a short introduction of the ES algorithms, the multimodal test functions will be introduced in Section III. In Section IV, the frozen noise model will be developed. Equations for the success probability will be derived in Section V. In Section VI, the scaling behavior of the population size will be investigated. Finally, a summary of the results and an outlook on future research will be given.

II. ES-ALGORITHMS

The basic $(\mu/\mu_I, \lambda)$ -ES algorithms investigated here consist of μ parents and λ offspring with truncation ratio $\vartheta := \mu/\lambda$. The subscript “ $m; \lambda$ ” denotes the selection of the $m = 1, \dots, \mu$ best individuals out of λ . The control of the strength σ of the isotropic Gaussian mutations used is done by either

Manuscript received 23 February 2024; revised 14 May 2024 and 18 June 2024; accepted 19 June 2024. Date of publication 27 June 2024; date of current version 8 October 2025. This work was supported by the Austrian Science Fund (FWF) under Grant P33702-N. This article was approved by Associate Editor A. Agapie. (Corresponding author: Lisa Schönerberger.)

The authors are with the Research Center Business Informatics, Vorarlberg University of Applied Sciences, 6850 Dornbirn, Austria (e-mail: lisa.schoenenberger@fhv.at; hans-georg.beyer@fhv.at).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TEVC.2024.3419931>, provided by the authors.

Digital Object Identifier 10.1109/TEVC.2024.3419931

Algorithm 1 $(\mu/\mu_l, \lambda)$ - σ SA Evolution Strategy

```

1: Initialize( $\mathbf{y}^{(0)}, \sigma^{(0)}, \sigma_{\text{stop}}, g = 0$ )
2: repeat
3:   for  $l = 1$  to  $\lambda$  do
4:      $\tilde{\sigma}_l = \sigma^{(g)} e^{\tau \mathcal{N}(0,1)}$  //mutate parental  $\sigma$ 
5:      $\tilde{\mathbf{y}}_l = \mathbf{y}^{(g)} + \tilde{\sigma}_l(\mathcal{N}(0,1), \dots, \mathcal{N}(0,1))$  //mutate  $\mathbf{y}$ 
6:      $\tilde{F}_l = F(\tilde{\mathbf{y}}_l)$  //evaluate offspring
7:   end for
8:   Sort Individuals  $\tilde{\mathbf{y}}$  Ascendingly w.r.t. Fitness  $\tilde{F}$ 
9:    $g = g + 1$ 
10:   $\mathbf{y}^{(g)} = \frac{1}{\mu} \sum_{m=1}^{\mu} \tilde{\mathbf{y}}_{m;\lambda}$  //recombine the  $\mu$  best  $\tilde{\mathbf{y}}$ 
11:   $\sigma^{(g)} = \frac{1}{\mu} \sum_{m=1}^{\mu} \tilde{\sigma}_{m;\lambda}$  //recombine the  $\mu$  best  $\tilde{\sigma}$ 
12: until  $\sigma^{(g)} < \sigma_{\text{stop}}$ 

```

Algorithm 2 $(\mu/\mu_l, \lambda)$ -CSA Evolution Strategy

```

1: Initialize( $\mathbf{y}^{(0)}, \sigma^{(0)}, \sigma_{\text{stop}}, \mathbf{s} = \mathbf{1}, g = 0$ )
2: repeat
3:   for  $l = 1$  to  $\lambda$  do
4:      $\tilde{\mathbf{z}}_l = (\mathcal{N}(0,1), \dots, \mathcal{N}(0,1))$  //search direction
5:      $\tilde{\mathbf{y}}_l = \mathbf{y}^{(g)} + \sigma^{(g)} \tilde{\mathbf{z}}_l$  //mutate  $\mathbf{y}$ 
6:      $\tilde{F}_l = F(\tilde{\mathbf{y}}_l)$  //evaluate offspring
7:   end for
8:   Sort Individuals  $\tilde{\mathbf{y}}$  Ascendingly w.r.t. Fitness  $\tilde{F}$ 
9:    $g = g + 1$ 
10:   $\mathbf{y}^{(g)} = \frac{1}{\mu} \sum_{m=1}^{\mu} \tilde{\mathbf{y}}_{m;\lambda}$  //recombine the  $\mu$  best  $\tilde{\mathbf{y}}$ 
11:   $\mathbf{z}^{(g)} = \frac{1}{\mu} \sum_{m=1}^{\mu} \tilde{\mathbf{z}}_{m;\lambda}$  //recombine the  $\mu$  best  $\tilde{\mathbf{z}}$ 
12:   $\mathbf{s} = (1 - c)\mathbf{s} + \sqrt{\mu c(2 - c)} \mathbf{z}^{(g)}$  //update  $\mathbf{s}$ -path
13:   $\sigma^{(g)} = \sigma^{(g-1)} \exp\left(\frac{\|\mathbf{s}\|^2 - N}{2DN}\right)$  //update  $\sigma$ , see [6, p.13]
14: until  $\sigma^{(g)} < \sigma_{\text{stop}}$ 

```

σ self-adaptation (σ SA), see Algorithm 1, or CSA, see Algorithm 2. For the σ SA-ES, the offspring mutation strengths are sampled from a log-normal distribution with learning parameter τ . The standard choice of the learning parameter is $\tau = 1/\sqrt{2N}$, which guarantees optimal performance on the sphere model [4]. A smaller τ yields a slower adaptation. This is an advantage in case of multimodal problems, which is one result of the investigations here. Therefore, a smaller learning rate of $\tau = 1/\sqrt{8N}$ will also be considered. For the CSA-ES, the standard choice for the cumulation time parameter c is $1/N$ and $1/\sqrt{N}$ [5], [6], where the latter results in faster convergence but a lower success probability P_s , which is also confirmed in the following sections.

III. MULTIMODAL TEST FUNCTIONS

The multimodal test functions considered have a high number of local optima: Rastrigin, Bohachevsky, and Ackley. The Rastrigin function $F_{\mathcal{R}}$ for an N -dimensional search vector $\mathbf{y} = (y_1, \dots, y_N)$ is given by

$$F_{\mathcal{R}}(\mathbf{y}) = \sum_{i=1}^N \left[y_i^2 + A(1 - \cos(\alpha y_i)) \right] \quad (1)$$

where the parameter A denotes the oscillation amplitude and α denotes the frequency. Unless otherwise stated, the parameters

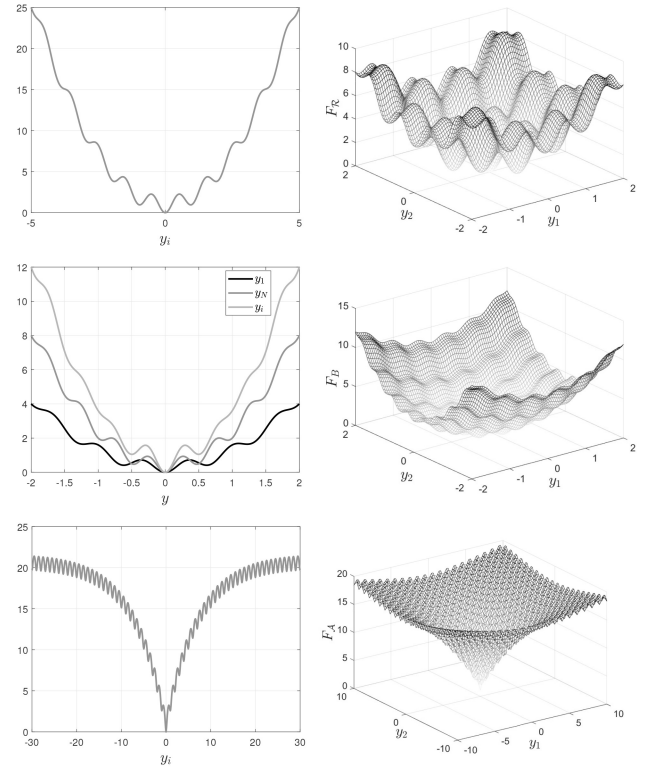


Fig. 1. Multimodal test functions for $N = 1$ and $N = 2$. From top to bottom: Rastrigin (1), Bohachevsky (2), and Ackley (3).

A and α have in all experiments the values $A = 1$ and $\alpha = 2\pi$. The global optimizer is located at $\hat{\mathbf{y}} = \mathbf{0}$ and is surrounded by local minima. The number of local minima is $\kappa^N - 1$ where κ increases with A and α (e.g., for $\alpha = 2\pi$, $A = 1$: $\kappa = 7$ and for $\alpha = 2\pi$, $A = 10$: $\kappa = 63$). The top plots of Fig. 1 show the Rastrigin function with standard parameters for $N = 1$ and $N = 2$, respectively.

The Bohachevsky function is given by

$$F_{\mathcal{B}}(\mathbf{y}) = \sum_{i=1}^{N-1} \left[y_i^2 + 2y_{i+1}^2 - B_1 \cos(\beta_1 y_i) - B_2 \cos(\beta_2 y_{i+1}) + B_1 + B_2 \right]. \quad (2)$$

The parameters B_1, B_2, β_1 , and β_2 control the size and frequency of the oscillations of the function. Larger values of B_1 and B_2 result in larger oscillations. Larger values of β_1 and β_2 increase the frequency and therefore the number of local minima. The standard values of these parameters are $B_1 = 0.3$, $B_2 = 0.4$, $\beta_1 = 3\pi$, and $\beta_2 = 4\pi$. Unless otherwise stated, these standard values will be used for all experiments. Compared to the Rastrigin function, the Bohachevsky function is not equally shaped in each dimension, as visualized in the middle-left plot of Fig. 1. It represents single coordinate contributions of the Bohachevsky function.

The Ackley function is given by

$$F_{\mathcal{A}}(\mathbf{y}) = C_1 - C_1 e^{-C_2 \sqrt{\frac{1}{N} \sum_{i=1}^N y_i^2}} + e - e^{\frac{1}{N} \sum_{i=1}^N \cos(\gamma y_i)} \quad (3)$$

and is represented by the bottom plots of Fig. 1 for $N = 1$ and $N = 2$. The function describes a funnel-shaped surface

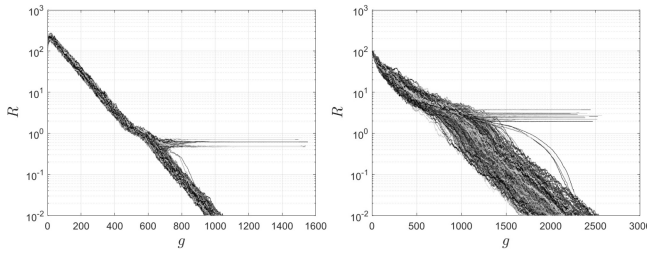


Fig. 2. R -dynamics for 500 $(\mu/\mu_I, 2\mu)$ - σ SA-ES runs with $\tau = 1/\sqrt{8N}$ for $N = 100$. Left plot: Bohachevsky landscape with $\mu = 50$. Right plot: Ackley landscape with $\mu = 5$. The success probabilities are $P_s = 0.88$ and $P_s = 0.9$, respectively. For each run, the ES was initialized randomly with distance $R^{(0)} = 100$ to the global optimizer.

with an infinite number of local minima. Larger values of C_1 lead to larger horizontal asymptotes. The larger C_2 , the steeper is the slope around the origin. The frequency is determined by γ . The standard values for the Ackley parameters are $C_1 = 20$, $C_2 = 0.2$, and $\gamma = 2\pi$. Unless otherwise stated, these parameters will be used in all subsequent experiments.

Fig. 2 shows the dynamics of the distance of the parental centroid to the global optimizer, i.e., $R(g) := \|\mathbf{y}^{(g)}\|$ of 500 independent σ SA-ES runs, Algorithm 1, in the Bohachevsky landscape (left plot) and in the Ackley landscape (right plot). Most runs converge to the global optimizer, but a certain percentage of runs is getting trapped in local minima. The percentage of runs that approach the global optimizer (i.e., the successful runs) is referred to as *success probability* P_s . The mean value dynamics of the successful runs in Fig. 2 are displayed in Fig. 3 and symbolized by angular brackets. The gray solid lines represent the distance R to the optimizer, mutation strength σ and its normalization

$$\sigma^* := \sigma \frac{N}{R}. \quad (4)$$

The values of σ^* are represented on the right y-axis. Having a closer look at the dynamics, one sees that the evolution process can be divided into different phases. At the beginning the initial parental centroid is initialized far enough from the global optimizer. Self-adaptation decreases the normalized σ^* for the Bohachevsky function, where the initial σ was chosen too large. For the Ackley function, on the other hand, the initial σ was (intentionally) chosen too small. Again σ^* reaches a nearly constant level that depends on the general structure of the function, which is further investigated in Section IV. Getting closer to the global optimizer, one observes a certain decrease of σ^* indicating the influence of the *local attractors* rewarding local fitness gain due to smaller mutation steps. That is, the ES reduces the normalized mutation strength σ^* . Finally, in the third phase, the ES is either trapped in a local attractor or it has hit the *global attractor* \mathcal{A}_{ES} , to be introduced in Section V-A. In the case of successful runs σ^* reaches, the constant level again. The conditions under which the global attractor is reached and the success probability P_s by which an ES approaches the global optimizer will be determined in Section V-B.

The solid black line in Fig. 3, shows the standard deviation of the residual part (to be defined in the following), which was

determined by experiments. The theoretical calculations for the standard deviation of the residual part are represented by the dashed-dotted black lines, with the horizontal lines indicating the asymptotic behavior. This will be further investigated in Section IV-C. Similar graphs can be obtained for the Rastrigin landscape. The CSA-ES, Algorithm 2, exhibits similar behaviors, see supplement material.

IV. FROZEN NOISE MODEL

A. Frozen Noise Model in a Nutshell

A theoretical analysis of the behavior of ES is usually based on progress rate analysis. This approach has been done for the Rastrigin function in [7], [8], and [9]. However, this approach is very involved and up to now restricted to the analysis of the expected value behavior. In this article, we use an alternative approach first proposed for the Rastrigin function in [3].

The evolution of the ES is interpreted as a walk through a frozen noise landscape. The derivation of the frozen noise model comprises two steps.

First, the function $F(\mathbf{y})$ to be optimized is subdivided into a global part $G(\mathbf{y})$ and a residual part $C(\mathbf{y})$ where the latter contains the nonconvex nonlinearities of F

$$F(\mathbf{y}) =: G(\mathbf{y}) + C(\mathbf{y}). \quad (5)$$

In this article, the global part is assumed to be expressible as a sphere model, i.e., $G(\mathbf{y}) = G(\|\mathbf{y}\|) = G(R)$.

Second, the residual part $C(\mathbf{y})$ is replaced by a noise term the distribution of which and its statistical parameters are to be determined. In order to further apply the noisy sphere theory [10], a Gaussian noise distribution will be used as approximation.

Determining the noise distribution is not arbitrary, but must reflect the real dynamics of the ES run. To this end it is important to recall how the ES explores the search space by a *restricted random walk* [11]. The ES mutations in line 5 of Algorithms 1 and 2 sample the fitness landscape in an unbiased manner. The direction of the evolution, however, is given by the selection. From this point of view, one may regard the evolution process as a restricted random walk because of the selection. Furthermore, each mutation and also the changes of the recombined parents (line 10 in Algorithms 1 and 2) do have a certain contribution toward the optimizer, but also contributions perpendicular to it. The length of this perpendicular part is on average by a factor of \sqrt{N} larger than the contribution toward the optimizer [11]. Furthermore, after the ES has reached the steady state, each component of \mathbf{y} is normally distributed with mean zero [10]. As a result, the parental centroids circle around the global optimizer (provided that the mutation strength is sufficiently large). This may be regarded as a sampling process, where the distance R to the global optimizer decreases only slightly from one generation to the next. If one neglects this small R decrease, i.e., considers the mutation process under the assumption of $R = \|\tilde{\mathbf{y}}\| = \text{const}$. Then, the fitness distribution of the offspring and therefore its standard deviation (below denoted as σ_{ES} to indicate that the noise is generated by the mutations of the ES) is solely determined by the residual part $C(\tilde{\mathbf{y}})$ of F . Thus,

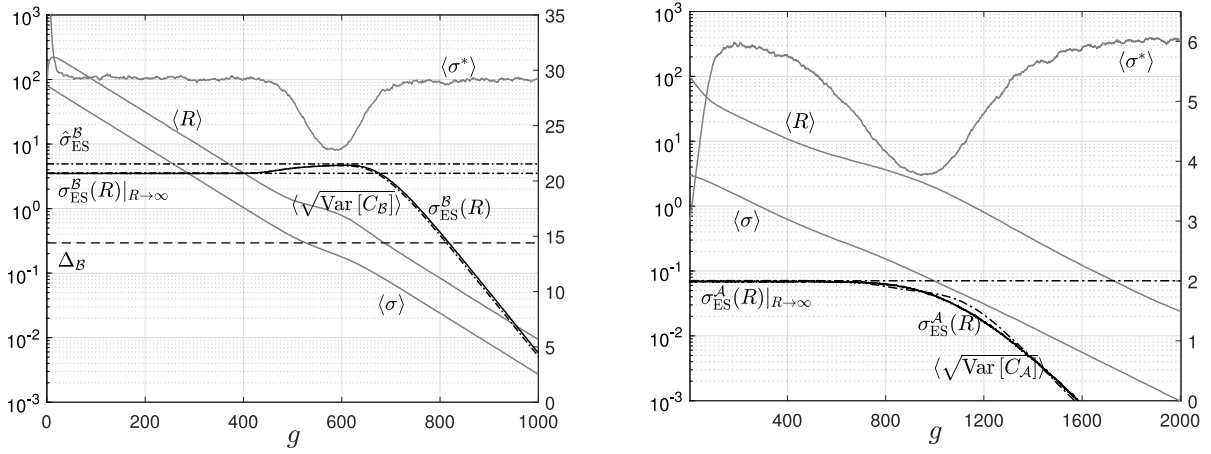


Fig. 3. Mean value dynamics derived from the *successful* runs displayed in Fig. 2. Left plot: Bohachevsky landscape with $\mu = 50$. Right plot: Ackley landscape with $\mu = 5$. In addition to the R dynamics, the mutation strength σ and its normalization σ^* (right y-axis) are displayed. $\sigma_{ES}(R)$, (15), (16), $\sigma_{ES}(R)|_{R \rightarrow \infty}$, (25), (33), $\hat{\sigma}_{ES}^B$, (26), and $\langle \sqrt{\text{Var}[C]} \rangle$ are discussed in Section IV-C. The distance Δ_B (40) is introduced in Section V-A.

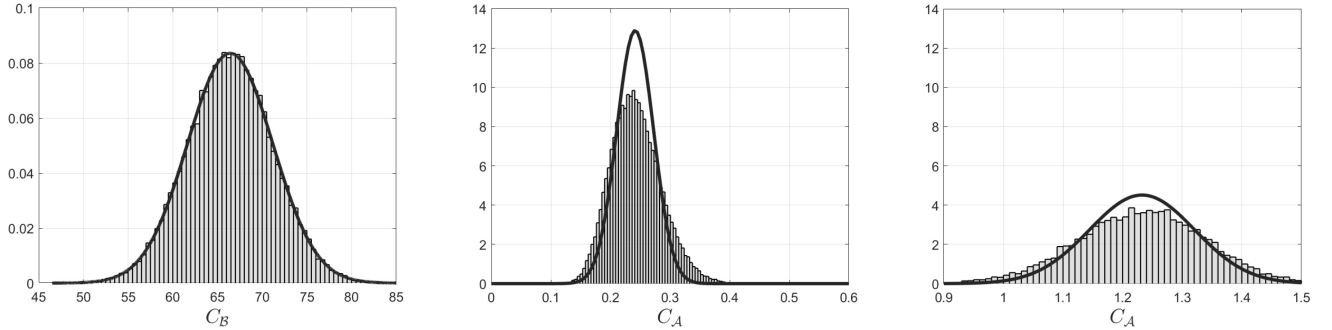


Fig. 4. Histogram of residual parts for 300 (50/50), 100-CSA-ES runs with $c = 1/\sqrt{N}$ and $N = 100$ in the Bohachevsky landscape at $R_{st}^B = 1.02$, (36) (left plot) and in the Ackley landscape at $R = 0.22$, (middle plot) and $R = 0.27$ (right plot). Bold solid lines show the pdf of normally distributed variates.

the ES sees basically a sphere model disturbed by distance depending noise. This defines the *frozen noise model*

$$\tilde{F}(\mathbf{y}) := G(R) + \sigma_{ES}(R)\mathcal{N}(0, 1), \quad R = \|\mathbf{y}\| \quad (6)$$

consisting of the global sphere model part $G(\mathbf{y}) = G(R)$ and a normally distributed perturbation part where its conditional variance is defined as

$$\sigma_{ES}^2(R) := \text{Var}[C(\mathbf{y} + \sigma \mathbf{z}) | \|\mathbf{y}\| = R], \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (7)$$

The assumption of a normal distribution in (6) must be generally regarded as an approximation. However, there is empirical evidence that this is a reasonable assumption (see, e.g., Fig. 4). In the following, $G(R)$ and $\sigma_{ES}(R)$ will be determined for the different test functions.

B. Building the Models

As for the Rastrigin function (1) the decomposition in terms of (5) reads

$$G_{\mathcal{R}}(\mathbf{y}) := \sum_{i=1}^N y_i^2 = R^2 \quad (8)$$

$$C_{\mathcal{R}}(\mathbf{y}) := NA - A \sum_{i=1}^N \cos(\alpha y_i) \quad (9)$$

where R is the distance to the global optimizer. The residual $C_{\mathcal{R}}(\mathbf{y})$ describes the oscillations of the Rastrigin function. It holds $C_{\mathcal{R}}(\mathbf{y}) \in [0, 2NA]$.

The decomposition of the Bohachevsky function (2) reads

$$G_{\mathcal{B}}(\mathbf{y}) := y_1^2 + 3 \sum_{i=2}^{N-1} y_i^2 + 2y_N^2 = 3R^2 - 2y_1^2 - y_N^2 \quad (10)$$

$$C_{\mathcal{B}}(\mathbf{y}) := (N-1)(B_1 + B_2) - \sum_{i=1}^{N-1} [B_1 \cos(\beta_1 y_i) + B_2 \cos(\beta_2 y_{i+1})]. \quad (11)$$

Note that the global part describes an hyperellipsoid. For $N \rightarrow \infty$, one can neglect the distortion of the 1st and last y -component to get a sphere model that will be used as an approximation for smaller N (see also Fig. 5, left plot for an exemplary comparison). For the residual part, it holds $C_{\mathcal{B}}(\mathbf{y}) \in [0, 2(N-1)(B_1 + B_2)]$.

The Ackley function (3) yields

$$G_{\mathcal{A}}(\mathbf{y}) := C_1 - C_1 \exp\left(-C_2 \sqrt{\frac{1}{N} \sum_{i=1}^N y_i^2}\right) \quad (12)$$

$$C_{\mathcal{A}}(\mathbf{y}) := e - \exp\left(\frac{1}{N} \sum_{i=1}^N \cos(\gamma y_i)\right). \quad (13)$$

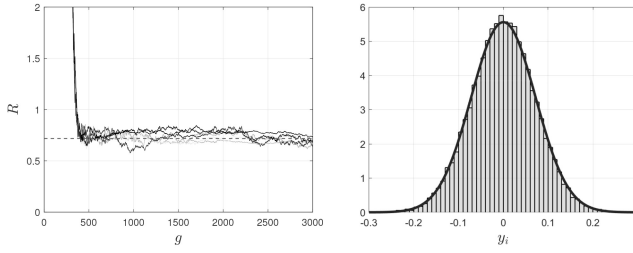


Fig. 5. Left plot: R -dynamics of the noisy hyperellipsoid model (10) and (11) (dark lines) and noisy sphere model (15) (light lines) using $\hat{\sigma}_{\text{ES}}^{\text{B}}$ from (26). Black dashed line shows R_{st}^{B} . Right plot: Histogram of all individual components y_i at distance R_{st}^{B} . Bold solid line shows the pdf of the $\mathcal{N}(0, R_{\text{st}}^2/N)$ variate. Experiments executed for the (100/100 $_f$, 200)-CSA-ES with $c = 1/\sqrt{N}$, $N = 100$, and $R_{\text{st}}^{\text{B}} = 0.72$, (36).

The global part is a nonquadratic sphere model with a funnel-shaped surface (see Fig. 1). For the residual part, one finds $C_{\mathcal{A}}(\mathbf{y}) \in [0, e - e^{-1}]$.

Using (6) together with the global parts (8), (10), and (12), respectively, one obtains the corresponding frozen noise models for the Rastrigin landscape

$$\tilde{F}_{\mathcal{R}}(\mathbf{y}) = R^2 + \sigma_{\text{ES}}^{\mathcal{R}}(R) \mathcal{N}(0, 1) \quad (14)$$

for the Bohachevsky landscape

$$\tilde{F}_{\mathcal{B}}(\mathbf{y}) = 3R^2 + \sigma_{\text{ES}}^{\mathcal{B}}(R) \mathcal{N}(0, 1) \quad (15)$$

and the Ackley landscape

$$\tilde{F}_{\mathcal{A}}(\mathbf{y}) = C_1 - C_1 e^{-C_2 \frac{R}{\sqrt{N}}} + \sigma_{\text{ES}}^{\mathcal{A}}(R) \mathcal{N}(0, 1) \quad (16)$$

where $R = \|\mathbf{y}\|$. The calculation of σ_{ES} and the assumption of a normal distribution are examined in the following.

C. Distribution of the Residual Parts

In the noisy sphere models (14)–(16), the assumption of Gaussian noise was made. For the Rastrigin and Bohachevsky function, this can be justified by considering the residual parts $C_{\mathcal{R}}(\mathbf{y})$, (9), and $C_{\mathcal{B}}(\mathbf{y})$, (11) as a sum of independent random variables $\cos(\alpha \tilde{y}_i)$ for which the central limit theorem of statistics holds in the limit $N \rightarrow \infty$. The assumption of a normal distribution is confirmed by the left plot in Fig. 4 where the histogram of the residual part (11) of the Bohachevsky landscape is shown. The histogram was obtained at the critical distance R_{st}^{B} (36) where the probability of getting trapped into a local attractor gets larger. The bold line shows the pdf of a normal distribution where the mean and variance were taken from the experimental data.

The calculation of the noise variances needed in (14)–(16) starts from (7) by first determining the variance of $C(\mathbf{y} + \sigma \mathbf{z})$ for a fixed parental centroid state \mathbf{y} and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\sigma_{\text{ES}}^2(\mathbf{y}) = \text{Var}[C(\mathbf{y} + \sigma \mathbf{z})]. \quad (17)$$

It is assumed that each offspring has the same mutation strength σ . This is exactly fulfilled for the CSA-ES (Algorithm 2, line 5), for the σ SA-ES it holds asymptotically

¹Here, we have used \tilde{y}_i to indicate the offspring components (see line 5 in Algorithms 1 and 2) that are responsible for the landscape sampling of the ES.

for $N \rightarrow \infty$. In a second step, averaging over the hypersphere $\|\mathbf{y}\| = R$ must be performed

$$\sigma_{\text{ES}}^2(R) = \int_{\|\mathbf{y}\|=R} \sigma_{\text{ES}}^2(\mathbf{y}) d^N y. \quad (18)$$

This can be done only for $N \leq 2$ in closed form, however, for $N \rightarrow \infty$ one can express (18) as an expected value of a normally distributed random vector [9]

$$\mathbf{y} \sim \mathcal{N}\left(\mathbf{0}, \frac{R^2}{N} \mathbf{I}\right) \quad (19)$$

$$\sigma_{\text{ES}}^2(R) = \mathbb{E}[\sigma_{\text{ES}}^2(\mathbf{y})]. \quad (20)$$

The variance and expected value calculations in (17) and (20) are straightforward but lengthy. They are detailed in supplement Section I. Only the calculation of (17) for the residual (9) of Rastrigin will be sketched. Using (9) in (17), one obtains

$$\begin{aligned} (\sigma_{\text{ES}}^{\mathcal{R}}(\mathbf{y}))^2 &= A^2 \text{Var}\left[\sum_{i=1}^N \cos(\alpha y_i + \alpha \sigma z_i)\right] \\ &= A^2 \sum_{i=1}^N \text{Var}[\cos(\alpha y_i + \alpha \sigma z_i)]. \end{aligned} \quad (21)$$

Here, the statistical independences of the $z_i \sim \mathcal{N}(0, 1)$ have been taken into account. Since the argument in the cosine can be interpreted as a random variate $x \sim \mathcal{N}(\alpha y_i, (\alpha \sigma)^2)$ one can show (see supplement)

$$\begin{aligned} \text{Var}[\cos(\alpha y_i + \alpha \sigma z_i)] \\ = \frac{1}{2} \left(1 - e^{-(\alpha \sigma)^2}\right) \left(1 - e^{-(\alpha \sigma)^2} \cos(2\alpha y_i)\right). \end{aligned} \quad (22)$$

Taking further $\mathbb{E}[\cos(2\alpha y_i)] = \exp(-[2(\alpha R)^2]/N)$ into account and substituting this into (22) and in turn into (21), one finally obtains for (20)

$$\left(\sigma_{\text{ES}}^{\mathcal{R}}(R)\right)^2 = A^2 \frac{N}{2} \left(1 - e^{-(\alpha \sigma)^2}\right) \left(1 - e^{-(\alpha \sigma)^2} e^{-\frac{2(\alpha R)^2}{N}}\right). \quad (23)$$

This result agrees well with runs of the CSA-ES and shows only small deviations for the σ SA-ES (see Fig. 1 in the supplement). The maximum noise is obtained for $R \rightarrow \infty$ since $\sigma = \sigma^* R/N$ one finds using (23)

$$\hat{\sigma}_{\text{ES}}^{\mathcal{R}} := \max_R [\sigma_{\text{ES}}^{\mathcal{R}}(R)] = \sigma_{\text{ES}}^{\mathcal{R}}(R)|_{R \rightarrow \infty} = A \sqrt{\frac{N}{2}}. \quad (24)$$

It will serve as an upper bound in subsequent calculations.

For the Bohachevsky function, the calculation of $\sigma_{\text{ES}}^{\mathcal{B}}(R)$ using (11) is similar, but there are additional covariance terms. The resulting rather large variance expression can be found in the supplement. It holds

$$\sigma_{\text{ES}}^{\mathcal{B}}(R)|_{R \rightarrow \infty} = \sqrt{\frac{1}{2} (N-1) (B_1^2 + B_2^2)} \quad (25)$$

and the maximum value of $\sigma_{\text{ES}}^{\mathcal{B}}(R)$ can be approximated as

$$\begin{aligned} \hat{\sigma}_{\text{ES}}^{\mathcal{B}} &:= \max_R [\sigma_{\text{ES}}^{\mathcal{B}}(R)] \\ &\approx \sqrt{\frac{1}{2} (N-1) (B_1 + B_2)^2 - B_1 B_2} \end{aligned} \quad (26)$$

(see supplement Section I-B). These values are represented by the dashed-dotted lines in the left plot of Fig. 3.

In the case of the Ackley function, one observes that the cosine sum appears in the exponent of (13). It is shown in supplement Section I-C that the residual part of the Ackley function becomes asymptotically a log-normally distributed random variate

$$C_{\mathcal{A}}(\mathbf{y}) \approx e - e^{\mathcal{N}(\mu_{\mathcal{A}}, \sigma_{\mathcal{A}})}. \quad (27)$$

This log-normal distribution can be approximated by a normally distributed random variable with mean and variance

$$\mathbb{E}[C_{\mathcal{A}}] = e - e^{\mu_{\mathcal{A}} + \frac{1}{2}\sigma_{\mathcal{A}}^2} \quad (28)$$

$$\text{Var}[C_{\mathcal{A}}] = e^{2\mu_{\mathcal{A}} + \sigma_{\mathcal{A}}^2} (e^{\sigma_{\mathcal{A}}^2} - 1) \quad (29)$$

where

$$\mu_{\mathcal{A}} = e^{-\frac{1}{2}(\sigma\gamma)^2} e^{-\frac{1}{2}\frac{(\gamma R)^2}{N}} \quad (30)$$

$$\sigma_{\mathcal{A}}^2 = \frac{1}{2N} \left(1 - e^{-(\sigma\gamma)^2}\right) \left(1 - e^{-(\sigma\gamma)^2} e^{-\frac{1}{2}\frac{(\gamma R)^2}{N}}\right). \quad (31)$$

Therefore

$$\sigma_{\text{ES}}^{\mathcal{A}}(R) = \sqrt{\text{Var}[C_{\mathcal{A}}(R)]} = \sqrt{e^{2\mu_{\mathcal{A}} + \sigma_{\mathcal{A}}^2} (e^{\sigma_{\mathcal{A}}^2} - 1)} \quad (32)$$

(see supplement Section I-C). The middle and right plot of Fig. 4 depict the experimentally determined distribution of the residual part at two different distances. The bold solid line shows the pdf of a normal distribution with mean and variance predicted by (28) and (29). The values of σ used for the calculation of $\mu_{\mathcal{A}}$ and $\sigma_{\mathcal{A}}$ were determined by the experimental values. The pdf of the corresponding log-normal distribution is not displayed since the differences to the normal distribution are nearly not visible. Similar to the other test functions the standard deviation of the residual part $C_{\mathcal{A}}$ approaches a maximum saturation value as $R \rightarrow \infty$ as can be seen in the rightmost graph of Fig. 3. Calculating the limit of (30), one sees that $\mu_{\mathcal{A}} \rightarrow 0$ and for (31) $\sigma_{\mathcal{A}}^2 \rightarrow (1/2N)$. Thus, one gets from (32) for $N \rightarrow \infty$ the saturation value

$$\hat{\sigma}_{\text{ES}}^{\mathcal{A}} := \sigma_{\text{ES}}^{\mathcal{A}}(R)|_{R \rightarrow \infty} = \frac{1}{\sqrt{2N}} \quad (33)$$

after a short calculation. This saturation value is depicted by the dashed-dotted line in the right plot of Fig. 3.

D. Steady-State in the Noise Model

It is known from the noisy sphere model theory that an ES optimizing an N -dimensional noisy sphere, $F_{\mathcal{S}} = aR^2 + \sigma_{\text{ES}}\mathcal{N}(0, 1)$, with constant noise strength σ_{ES} reaches a steady-state R -distribution with $R_{\text{st}} := \mathbb{E}[R] \neq 0$ where

$$R_{\text{st}} \simeq \sqrt{\frac{\sigma_{\text{ES}}N}{4a\mu c_{\mu/\mu, \lambda}}} \quad (34)$$

see [4], [12]. It also holds for the noisy sphere [10] that under steady-state conditions, each component of \mathbf{y} is normally distributed with

$$y_i = (\mathbf{y})_i \sim \mathcal{N}\left(0, \frac{R_{\text{st}}^2}{N}\right). \quad (35)$$

As already shown in [3], this also holds approximately for the parental distribution of the ES on Rastrigin as long as the ES is not trapped in one of the local attractors. In the case of the Bohachevsky landscape, the global part is a hyperellipsoid. As already state, one can assume that for large N the first and last component has only limited influence on the steady-state distribution. Neglecting the first and last component, one gets the noisy sphere model (15). Using (34), the expected value of the steady-state becomes

$$R_{\text{st}}^{\mathcal{B}} \simeq \sqrt{\frac{\hat{\sigma}_{\text{ES}}^{\mathcal{B}}N}{12\mu c_{\mu/\mu, \lambda}}}. \quad (36)$$

This is illustrated in Fig. 5 for dimension $N = 100$, where the left plot shows that the value of $R_{\text{st}}^{\mathcal{B}}$ predicts the residual distance of the experimental data well. Equation (35) also holds approximately for the Bohachevsky function, as illustrated in the right plot of Fig. 5. The histogram shows the distribution of all components lumped together. The variance of the experimental values is 0.0052 as predicted. Note that the distributions of the single first and single last component have larger variances which can be neglected. The influence of this simplification is investigated in supplement Section III-B.

V. SUCCESS PROBABILITY MODEL

A. Global Attractor Region

Using the models (14)–(16), the evolution of the ES in the multimodal landscape can be interpreted as a noisy minimization problem. Recall that the ES reaches the vicinity of the global optimizer up to a distance R_{st} provided that the noise variance does not change. However, if this distance is close enough to the global optimizer the noise strength declines (see Fig. 3) and a large success probability P_s can be expected. To specify what is close enough, the following definition is introduced.

The attractor region is defined as the region where the negative gradient flow goes toward the optimizer. In [3], it was already shown for the Rastrigin function that the attractor region is a hypercube

$$\mathcal{A}_{\text{ES}}^{\mathcal{R}} := [-\Delta_{\mathcal{R}}, \Delta_{\mathcal{R}}]^N \quad (37)$$

where $\Delta_{\mathcal{R}}$ is the distance from the global optimizer to the nearest stationary points. It holds

$$\Delta_{\mathcal{R}} \simeq \frac{A\alpha\pi}{A\alpha^2 - 2}. \quad (38)$$

This is visualized for $N = 2$ in the left plot of Fig. 6. For the Bohachevsky function, the saddle point $\Delta_{\mathcal{B}}^{(i)}$ next to the global optimizer in coordinate direction $i \notin \{1, N\}$ is given by the 2nd zero of the derivative of (2), i.e.,

$$0 = 6y_i + B_1\beta_1 \sin(\beta_1 y_i) + B_2\beta_2 \sin(\beta_2 y_i) \Leftrightarrow y_i = \Delta_{\mathcal{B}}^{(i)}. \quad (39)$$

This nonlinear equation can be approximately solved by expanding the expressions $B \sin(\beta y_i)$ into a Taylor series at $x = \pi/\beta$, i.e., $-\sin(\beta y_i) = -\sin(\beta x) - \beta \cos(\beta x)(y_i - x) + \mathcal{O}(y_i)^2 = \beta y_i - \pi + \mathcal{O}(y_i)^2$ and it follows after a short calculation:

$$\Delta_{\mathcal{B}} := \Delta_{\mathcal{B}}^{(i)} = \frac{\pi(B_1\beta_1 + B_2\beta_2)}{B_1\beta_1^2 + B_2\beta_2^2 - 6}. \quad (40)$$

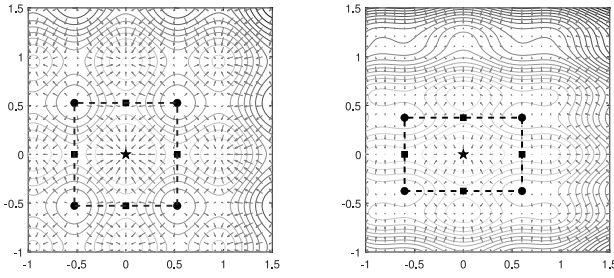


Fig. 6. Global attractor region of the Rastrigin (left) and Bohachevsky (right) function with standard parameters for $N = 2$. Arrows show the negative gradient flow. The star shows the global optimizer, squares the nearest stationary points, and circles the farthest stationary points (magnified pictures can be found in supplement Section II-A Fig. 4).

The estimation of the distances to the nearest stationary point in the first and N th coordinate directions is similar, i.e.,

$$\Delta_B^{(1)} = \frac{B_1 \beta_1 \pi}{B_1 \beta_1^2 - 2}, \quad \Delta_B^{(N)} = \frac{B_2 \beta_2 \pi}{B_2 \beta_2^2 - 4} \quad (41)$$

and the attractor region is a parallelepiped

$$\mathcal{A}_{\text{ES}}^B := \left[-\Delta_B^{(1)}, \Delta_B^{(1)} \right] \times \left[-\Delta_B, \Delta_B \right]^{N-2} \times \left[-\Delta_B^{(N)}, \Delta_B^{(N)} \right]. \quad (42)$$

This is visualized in the right plot of Fig. 6.

For the Ackley function, it is not possible to give an exact equation for the attractor region. Only the stationary points that lie on the coordinate axis can be determined analytically. For the stationary points, that do not lie on the coordinate axes an analytical calculation is not feasible. In supplement Section II-B, the stationary points on the axes are calculated and it is shown that the attractor region is bounded. The Investigations in Section V-E will show that the frozen noise model is not directly applicable to calculate the success probability. Therefore, a different approach is used which is independent of the attractor size.

To ensure convergence to the global optimizer, the parental centroid should be located in the global attractor region, i.e., $\mathbf{y} \in \mathcal{A}_{\text{ES}}$. Although the attractor regions $\mathcal{A}_{\text{ES}}^R$ (37) or $\mathcal{A}_{\text{ES}}^B$ (42) for the Rastrigin and Bohachevsky function, respectively, can be specified exactly, the following considerations must be taken into account for further investigations: A gradient strategy, starting inside the attractor region is always successful. For an ES, this is not necessarily the case. Especially when the parental centroid is located in a corner, the probability to produce a better offspring is about 2^{-N} . On the other hand, when the parental centroid is located outside the attractor but in the vicinity of an edge, the probability of producing an offspring inside the attractor region is approximately² $1/2$. This leads to the introduction of a correction term ε that depends on the strategy-specific parameters, such as the current mutation strength σ , truncation ratio $\vartheta = \mu/\lambda$, learning parameter τ , or cumulation constant c , respectively.

²Suppose σ is chosen to be sufficiently small to avoid jumping over the attractor region, and sufficiently large to overcome the distance to the attractor region.

The attractor region containing the correction term is defined as

$$\mathcal{A}_{\text{ES}}^R(\varepsilon) := [-\Delta_R - \varepsilon, \Delta_R + \varepsilon]^N \quad (43)$$

for the Rastrigin function and

$$\mathcal{A}_{\text{ES}}^B(\varepsilon) := \left[-\Delta_B^{(1)} - \varepsilon, \Delta_B^{(1)} + \varepsilon \right] \times \left[-\Delta_B - \varepsilon, \Delta_B + \varepsilon \right]^{N-2} \times \left[-\Delta_B^{(N)} - \varepsilon, \Delta_B^{(N)} + \varepsilon \right] \quad (44)$$

for the Bohachevsky function, respectively. By introducing ε , the attractor model can be used for any ES version.

B. Estimating the Success Probability

As illustrated in the figures of Section III the noise strength σ_{ES} is bounded. Therefore, it can be deduced from (34) and (36) that the probability of reaching the attractor mainly depends on the choice of a sufficiently large population size μ (assuming $\vartheta = \text{const.}$). If the population size is too small, then R_{st} is too large and there is no intersection between the attractor region and the sphere with radius R_{st} . In this case, one can assume that global convergence is virtually impossible. On the other hand, if μ is very large such that the R_{st} -sphere lies completely inside the attractor region, every ES reaches the attractor that leads to a success probability near one. For the cases, where the attractor region intersects the R_{st} -sphere positive progress is possible with a certain success probability P_s .

As has been shown in Section IV-D (35), reaching the steady-state distance R_{st} the single parental components in the Rastrigin and Bohachevsky landscape are approximately normally distributed with zero mean and standard deviation

$$\sigma_{\text{st}} := \frac{R_{\text{st}}}{\sqrt{N}}. \quad (45)$$

For the Rastrigin and Bohachevsky function, the attractor region containing the correction term ε can be described by (43) and (44). Taking this into account, one can derive an equation for the success probability, i.e., the probability that the parental centroid is located in the attractor region under the assumption that the single components are normally distributed with mean zero and standard deviation σ_{st} . Using the independence of the parental centroid components in the steady-state, it follows for the Bohachevsky function:

$$\begin{aligned} P_s^B &= \Pr[\mathbf{y} \in \mathcal{A}_{\text{ES}}^B] = \Pr\left[\left(-\Delta_B^{(1)} - \varepsilon \leq y_1 \leq \Delta_B^{(1)} + \varepsilon \right) \right. \\ &\quad \wedge \left(-\Delta_B - \varepsilon \leq y_2 \leq \Delta_B + \varepsilon \right) \wedge \dots \\ &\quad \left. \wedge \left(-\Delta_B^{(N)} - \varepsilon \leq y_N \leq \Delta_B^{(N)} + \varepsilon \right) \right] \\ &\approx \Pr[-\Delta_B - \varepsilon \leq y \leq \Delta_B + \varepsilon]^N. \end{aligned} \quad (46)$$

In the last approximation, the influence of the first and last component was neglected, the attractor region is assumed as a hypercube. This simplification has only a small impact on further calculations as it is shown in the supplement Section III-B. Using the fact that the single components are normally distributed, one gets for a single component

$$\begin{aligned}
& \Pr[-\Delta_{\mathcal{B}} - \varepsilon \leq y \leq \Delta_{\mathcal{B}} + \varepsilon] \\
&= \Pr\left[-\frac{\Delta_{\mathcal{B}} + \varepsilon}{\sigma_{\text{st}}} \leq z \leq \frac{\Delta_{\mathcal{B}} + \varepsilon}{\sigma_{\text{st}}}\right] \\
&= \Phi\left(\frac{\Delta_{\mathcal{B}} + \varepsilon}{\sigma_{\text{st}}}\right) - \Phi\left(-\frac{\Delta_{\mathcal{B}} + \varepsilon}{\sigma_{\text{st}}}\right) \quad (47)
\end{aligned}$$

where $\Phi(z)$ is the cdf of the standard normal variate $\mathcal{N}(0, 1)$ and σ_{st} is given by (45). Thus, one gets for the success probability (46)

$$P_s^{\mathcal{B}} = \left[2\Phi\left(\frac{\Delta_{\mathcal{B}} + \varepsilon}{\sigma_{\text{st}}}\right) - 1\right]^N \quad (48)$$

where σ_{st} depends on R_{st} and therefore on the noise strength $\sigma_{\text{ES}}^{\mathcal{B}}(R_{\text{st}})$. Considering the left plot of Fig. 3, the general course of the curve is that $\sigma_{\text{ES}}^{\mathcal{B}}(R)$ first stays constant for large R followed by a small increase and then starts to drop in the region where the distance R is of the order of $\Delta_{\mathcal{B}}$. That is, even if the ES enters the global attractor region $\mathcal{A}_{\text{ES}}^{\mathcal{B}}$, $\sigma_{\text{ES}}^{\mathcal{B}}(R)$ is still in the vicinity of its maximum value. Therefore, the maximum value $\hat{\sigma}_{\text{ES}}^{\mathcal{B}}$ (26) for the Bohachevsky function is used and it follows from (36) and (45):

$$\sigma_{\text{st}}^{\mathcal{B}} = \sqrt{\frac{\hat{\sigma}_{\text{ES}}^{\mathcal{B}}}{12\mu c_{\mu/\mu, \lambda}}} = \sqrt{\frac{\sqrt{\frac{1}{2}(N-1)(B_1+B_2)^2 - B_1B_2}}{12\mu c_{\mu/\mu, \lambda}}}. \quad (49)$$

As already mentioned in Section IV, since the global part of the Bohachevsky function is a hyperellipsoid, the standard deviation of the first and last component is larger. However, it can again be assumed that this effect is neglectable for large N , which is demonstrated in more detail in the supplement Section III-B. Inserting this into (48), one finally obtains the success probability in the Bohachevsky case

$$P_s^{\mathcal{B}} = \left[2\Phi\left(\sqrt{\frac{12\mu c_{\mu/\mu, \lambda}}{\hat{\sigma}_{\text{ES}}^{\mathcal{B}}}}(\Delta_{\mathcal{B}} + \varepsilon)\right) - 1\right]^N \quad (50)$$

with $\hat{\sigma}_{\text{ES}}^{\mathcal{B}} = \sqrt{(1/2)(N-1)(B_1+B_2)^2 - B_1B_2}$. The derivation of an equation for the Rastrigin function is almost analogous and has already been done in [3]. Using (24) and (34) with $a = 1$ the success probability of the Rastrigin function becomes

$$P_s^{\mathcal{R}} = \left[2\Phi\left(\sqrt{\frac{4\sqrt{2}\mu c_{\mu/\mu, \lambda}}{A\sqrt{N}}}(\Delta_{\mathcal{R}} + \varepsilon)\right) - 1\right]^N. \quad (51)$$

C. Comparison With Experiments

Figs. 7 and 8 evaluate the predictive quality of the success probability (50), and (51) using $\Delta_{\mathcal{R}}$ from (38), $\Delta_{\mathcal{B}}$ from (40) and $\varepsilon = 0$. The equations are represented by the solid lines and are compared with experimental values represented by the stars. Each data point was obtained by at least 1000 independent ES runs. The experiments were executed for the σ SA-ES and the CSA-ES. As expected, there are differences between the experimental data and the predictions. However, in most cases the general tendencies are well covered by the predictions from (51) and (50). The shape of the curves

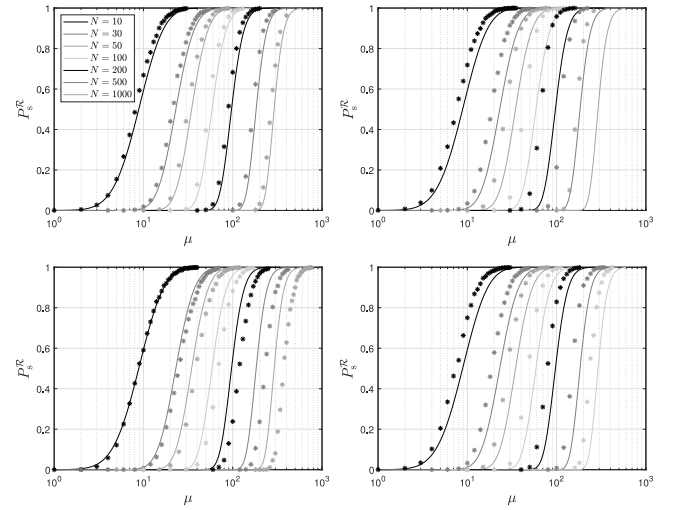


Fig. 7. Success $P_s^{\mathcal{R}}$ depending on population size μ for the Rastrigin landscape. Solid lines show $P_s^{\mathcal{R}}$ predicted by (51) with $\varepsilon = 0$ and $\vartheta = 0.5$. Experimental values are displayed by the stars. From top left to bottom right: CSA-ES with $c = 1/\sqrt{N}$, $c = 1/N$, σ SA-ES with $\tau = 1/\sqrt{2N}$, and $\tau = 1/\sqrt{8N}$.

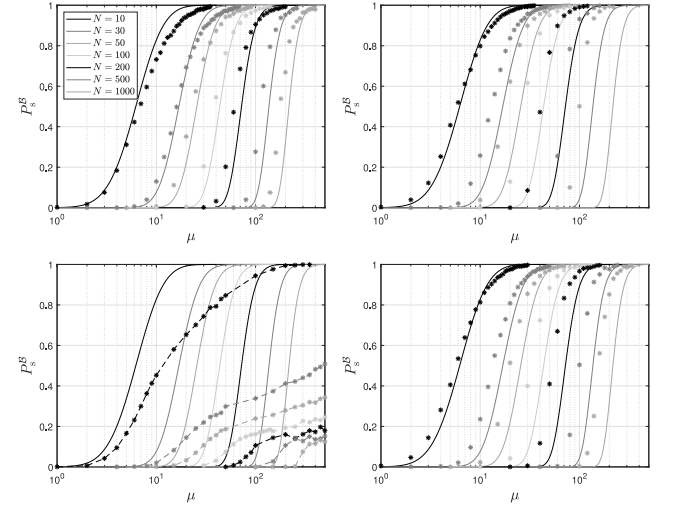


Fig. 8. Success $P_s^{\mathcal{B}}$ depending on population size μ for the Bohachevsky landscape. Solid lines show $P_s^{\mathcal{B}}$ predicted by (50) with $\varepsilon = 0$ and $\vartheta = 0.5$. Experimental values are displayed by the stars. From top left to bottom right: CSA-ES with $c = 1/\sqrt{N}$, $c = 1/N$, σ SA-ES with $\tau = 1/\sqrt{2N}$, and $\tau = 1/\sqrt{8N}$.

differs more in the case of the Bohachevsky function. The success probability increases slower than predicted. This is most evident in the case of the σ SA-ES with standard choice $\tau = 1/\sqrt{2N}$ where the prediction (50) is not usable.

For both Rastrigin and Bohachevsky, there are considerable differences between the strategies investigated here. The CSA-ES with cumulation constant $c = 1/N$ has a larger success probability than the CSA-ES with the larger cumulation constant $c = 1/\sqrt{N}$. This is also the case for the σ SA-ES where the smaller learning parameter $\tau = 1/\sqrt{8N}$ leads to larger success probabilities. For both strategy types, smaller strategy parameters lead to slower adaptation, which is advantageous in multimodal landscapes: The strength of the mutations remains stable longer and sufficiently large so that the probability of being caught in a local minimum is lower. These observations

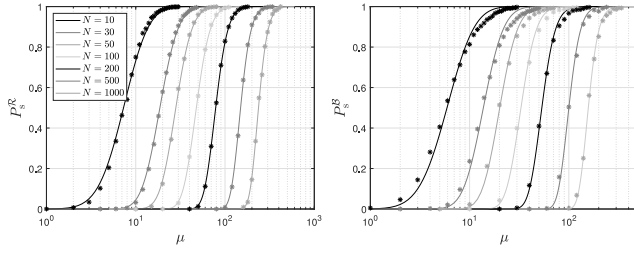


Fig. 9. Success P_s predicted by (50) and (51) with ε values according to Fig. 10 for σ SA-ES with $\tau = 1/\sqrt{8N}$ and $\vartheta = 0.5$ in the Rastrigin landscape (left plot) and in the Bohachevsky landscape (right plot).

do also hold for other truncation ratios, such as $\vartheta = 1/4$ (see supplement Section III-C).

D. Taking Strategy-Specific Properties into Account by ε

Due to the conceptional simplicity of the frozen noise model presented not all ES specific aspects can be modeled. That is, the actual size of the global attractor region depends on strategy-specific aspects, such as τ, c, ϑ , etc. To improve the prediction quality, the correction term $\varepsilon \neq 0$ is introduced in (50) and (51). The results for the σ SA-ES with $\tau = 1/\sqrt{8N}$ are presented in Fig. 9 where ε was chosen according to Fig. 10. The ε values were determined empirically by minimizing the sum of the squares of the differences between the equations and the experimental values. In the Rastrigin landscape, this leads to curves that do well agree with the real ES runs, as visible in the left plot of Fig. 9. For the Bohachevsky function (right plot), there are smaller deviations. Similar results are provided for the other strategy types in the supplement Section III-D.

E. Ackley Function

Since the size and shape of the global attractor region are not available it is not possible to derive an equation for the success probability as (50). Fig. 11 shows the experimental results for the success probability. Each data point was determined by at least 1000 ES runs. The dependence on N is much weaker than for the Rastrigin or Bohachevsky function. Especially for the CSA-ES with $c = 1/N$, the probabilities become only slightly smaller with increasing N . For the σ SA-ES with $\tau = 1/\sqrt{2N}$ and for large N the success probabilities are very low compared to the other strategy types. This is a result of the small initial mutation strength that prevents divergence but leads to local convergence in this case. All experiments were executed with moderate initial values $R^{(0)} = 10\sqrt{N}$ and $\sigma^{*(0)} = 3$. Compared to the other multimodal landscapes considered, the initial values in the Ackley landscape have a large influence on the success probability. The initial R distance must not be too large. This can be explained intuitively. Far outside the global part of the Ackley function the landscape is very flat, as visible in Fig. 1. In an almost flat landscape only small progress is possible due to the oscillations resulting in a large noise to signal ratio. Thus, no positive progress is possible and the ES diverges. This is also the case if the normalized mutation strength is too large. For the noisy sphere model, this has already been investigated. Therefore, an alternative

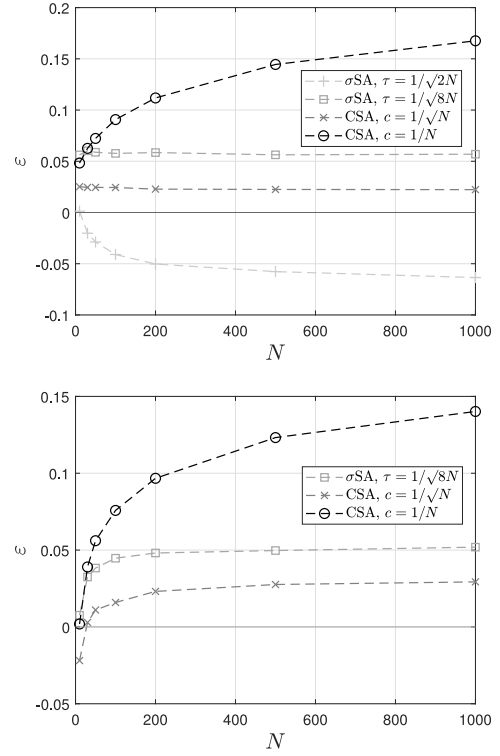


Fig. 10. Optimal ε that minimizes the deviations between (50) and (51) and the experimental values of Figs. 7 and 8. Top plot represents the results for the Rastrigin function and bottom plot for the Bohachevsky function, respectively.

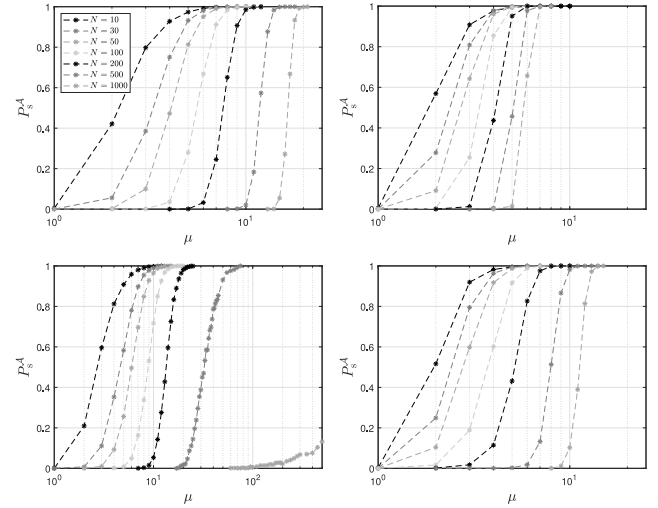


Fig. 11. Success P_s^A depending on population size μ in the Ackley landscape with $\vartheta = 0.5$. Experimental values are displayed by the stars. From top left to bottom right: CSA-ES with $c = 1/\sqrt{N}$, $c = 1/N$, σ SA-ES with $\tau = 1/\sqrt{2N}$ and $\tau = 1/\sqrt{8N}$.

model approach for the Ackley function can be derived from it, which is described as follows.

The maximal normalized mutation strength must fulfill the sufficient evolution condition [10]

$$(\sigma_{\text{ES}}^*)^2 + (\sigma^*)^2 \leq 4\mu^2 c_{\mu/\mu, \lambda}^2. \quad (52)$$

σ_{ES}^* denotes the normalized noise strength [13]

$$(\sigma_{\text{ES}}^A)^* = \sigma_{\text{ES}}^A(R) \frac{N}{RG'_A(R)} \quad (53)$$

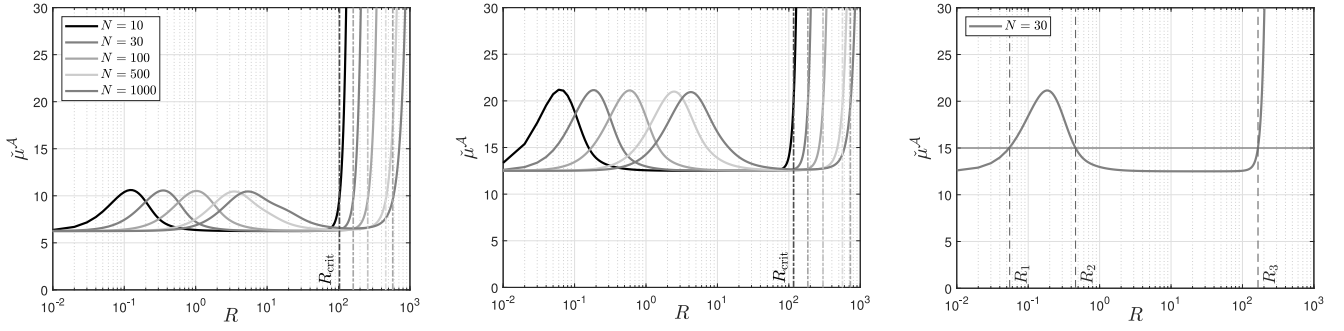


Fig. 12. Solid lines show $\check{\mu}^A$ from (54) with $\sigma^* = 10$ (left plot) and $\sigma^* = 20$ (middle and right plot). Vertical dashed-dotted lines in left and middle plot show R_{crit} (56).

where G'_A is the derivative w.r.t. R of the global part (12) of the Ackley function. Resolving the equation in (52) for μ using (53) taking (30)–(32) into account, one gets $\check{\mu}^A$ above which the ES can converge

$$\check{\mu}^A = \frac{\sqrt{[(\sigma_{\text{ES}}^A)^*]^2 + (\sigma^*)^2}}{2c_{\mu/\mu,\lambda}}. \quad (54)$$

This minimal parental population size is displayed in Fig. 12. The $\check{\mu}^A$ curves reveal a very interesting convergence behavior. For example, consider the $N = 30$, $\sigma^* = 20$ curve in the middle plot and draw a horizontal line at $\mu = 15$. One obtains three intersection points, $R_1 < R_2 < R_3$. These points are displayed in the right plot of Fig. 12. If the parental centroid of the ES is at a distance $R < R_1$ then the ES converges to the global optimizer. If R is located between R_1 and R_3 , the ES will converge to R_2 , i.e., it ends in a local attractor. If, however, $R > R_3$ the ES is expected to diverge. That is, the increase of the necessary population size $\check{\mu}^A$ beyond the R_3 value is almost exponential w.r.t. the R value: Using (52) with $\sigma^* = 0$ and (53) with (12) and (33), one finds

$$\mu \geq \frac{N}{2\sqrt{2}c_{\mu/\mu,\lambda}RC_1C_2} e^{\frac{C_2}{\sqrt{N}}R}. \quad (55)$$

Considering the curves for different search space dimensionalities N in Fig. 12 (left and middle), it becomes apparent that the value of the local maximum $\max(\check{\mu}^A)$ of the $\check{\mu}^A$ curves is rather insensitive w.r.t. N .³ This means that choosing a $\mu > \max(\check{\mu}^A)$ and an initial R being less than a critical initial R_{crit} , the global convergence is ensured independent of N . This is in contrast to the Rastrigin and Bohachevsky case.

To calculate the thresholds R_{crit} and σ_{crit}^* beyond which convergence becomes unattainable, insert the normalized noise strength (53) with (33) into the evolution condition (52). The resulting nonlinear equation can be solved approximately, yielding

$$R_{\text{crit}} = \frac{x_{\text{crit}}\sqrt{N}}{2C_2} \quad (56)$$

$$\sigma_{\text{crit}}^* = \sqrt{4\mu^2c_{\mu/\mu,\lambda}^2 - \frac{1}{2}\left(\frac{N}{RC_1C_2}\right)^2 e^{2C_2\frac{R}{\sqrt{N}}}} \quad (57)$$

³This N -independence also occurs when varying the Ackley parameters or ϑ .

where $x_{\text{crit}} \approx x^{(k)}$ is obtained recursively with

$$x^{(0)} = \ln\left(\frac{C_1(4\mu^2c_{\mu/\mu,\lambda}^2 - (\sigma^*)^2)}{2N}\right) \quad (58)$$

$$x^{(k)} = x^{(0)} + 2\ln(x^{(k-1)}). \quad (59)$$

Comprehensive calculations and experimental results are represented in supplement Section IV. Equation (56) is represented by the dashed lines in Fig. 12 (left and middle). It lies in the region where $\check{\mu}^A$ increases strongly.

The severe influence of the initial distance to the optimizer on the convergence behavior gives rise to the question how this value is usually chosen in test beds for the empirical evaluation of evolutionary algorithms. The start value for the CMA-ES in [1] is randomly initialized within the interval $[1, 30]^N$. The initial mutation strength σ is set to half of the initialization interval. Numerous other papers have adopted these starting values. R_{crit} (56) can fall within this initialization interval. For example, for $N = 100$, $\mu = 10$, and $\vartheta = 0.5$, it holds $R_{\text{crit}} \approx 270 < 300$, where 300 is the maximum distance in the initialization interval. Using Jensen's inequality, it follows for the norm of the vector with uniformly distributed components in $[a, b]$, that:

$$\begin{aligned} \mathbb{E}\left[\sqrt{\sum_{i=1}^N X_i^2}\right] &\leq \sqrt{\mathbb{E}\left[\sum_{i=1}^N X_i^2\right]} = \sqrt{\sum_{i=1}^N \mathbb{E}[X_i^2]} \\ &= \sqrt{N\frac{b^3 - a^3}{3(b-a)}} = \sqrt{\frac{N}{3}(a^2 + ab + b^2)}. \end{aligned} \quad (60)$$

The last step uses the second moment of the uniform distribution. Using the above mentioned values, the expected start distance is less than 176. For these parameters and for the CSA-ES with $c = 1/\sqrt{N}$, the runs already diverge from a starting distance of 150.

Swarm optimization algorithms often use a limited search space of $[-32.768, 32.768]^N$ and an initialization interval of $[16.384, 32.768]^N$ [14], $[-32.768, 32.768]^N$ [15], or $[-5, 5]^N$ [16]. Considering these values one might speculate that these initializations have been chosen within or at the edge of global convergence. That is, these choices are not suited to show the limitations of the algorithms. Therefore, test beds should be adopted to also show the limitations.

VI. POPULATION SIZING

A. Derivation of Parent Population Size

A main question of this article concerns the choice of μ and λ that guarantees convergence of the ES toward the global optimizer. As already mentioned in Section V-B the success probability mainly depends on the choice of the population size μ . A large μ leads to a small R_{st} and the R_{st} -sphere intersects the attractor region. Given a fixed truncation ratio ϑ , it suffices to derive a formula that predicts $\mu(P_s)$. It was shown in Section V-D that (50), (51) in conjunction with the correction term ε predict the experimental data well. Therefore, the equations can be used to derive a population sizing formula and to evaluate its scaling behavior. Solving (51) for μ under the assumption $c_{\mu/\mu,\lambda} \simeq f(\vartheta)$ [13, p. 249] yields for the Rastrigin function

$$\mu^{\mathcal{R}} \simeq \frac{\hat{\sigma}_{\text{ES}}^{\mathcal{R}}}{(\Delta_{\mathcal{R}} + \varepsilon)^2 4c_{\mu/\mu,\lambda}} \left[\Phi^{-1} \left(\frac{1}{2} + \frac{1}{2} P_s^{\frac{1}{N}} \right) \right]^2 \quad (61)$$

where Φ^{-1} is the quantile function of the standard normal distribution. For the Bohachevsky function, it follows from (50):

$$\mu^{\mathcal{B}} \simeq \frac{\hat{\sigma}_{\text{ES}}^{\mathcal{B}}}{(\Delta_{\mathcal{B}} + \varepsilon)^2 12c_{\mu/\mu,\lambda}} \left[\Phi^{-1} \left(\frac{1}{2} + \frac{1}{2} P_s^{\frac{1}{N}} \right) \right]^2. \quad (62)$$

Equations (61) and (62) with $\varepsilon = 0$ behave asymptotically like [3]

$$\mu = \mathcal{O}(\sqrt{N} \ln(N)). \quad (63)$$

As there is no analytical characterization of the success domain for the Ackley function, it is not possible to provide a μ -dependent equation for the success probability. However, (54) describes the minimal value of μ above which an ES can converge. Fig. 12 illustrates that there is a local maximum for $R < R_{\text{crit}}$. It follows that:

$$\mu^{\mathcal{A}} = \max_R [\tilde{\mu}^{\mathcal{A}}(R) | R < R_{\text{crit}}] \quad (64)$$

is an upper bound for the population size required to achieve a success rate of 1. Fig. 12 also shows that $\mu^{\mathcal{A}}$ strongly depends on σ^* . On the other hand, the value of the steady-state σ^* depends on μ and on strategy-specific parameters. Therefore, within the current model it is not possible to provide a general equation for $\mu^{\mathcal{A}}$. A value for σ^* must be assumed in order to determine $\mu^{\mathcal{A}}$.

B. Comparison With Experiments

Fig. 13 compares the predictions of the population size equations depending on N with experiments. At least 300 runs were executed to obtain the experimental data displayed by the markers. The solid lines represent the population sizing (61) and (62) with $\varepsilon = 0$ for $P_s = 50\%$ and for $P_s = 99\%$. While the theoretical predictions differ from the experimental values, (61) and (62) predict the general functional tendency well.

The deviations are due to the different values of ε , see Fig. 10. For the CSA-ES with $c = 1/N$, ε is positive and increases as N increases. This results in an overestimation of the population size. The scaling behavior is slower than

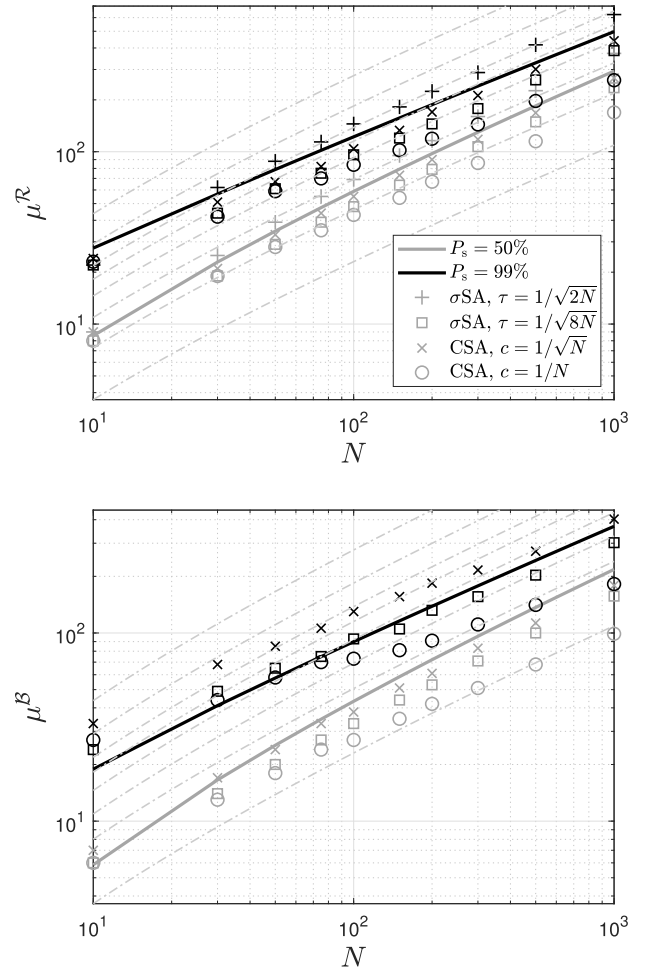


Fig. 13. Population size for $P_s = 50\%$ (lighter curve and markers) and $P_s = 99\%$ (darker curve and markers). Upper plot shows the results for Rastrigin and lower plot for Bohachevsky, both with $\vartheta = 0.5$. Bold lines show (61) and (62) and symbols represent the experimental results. Gray dashed-dotted lines show functions $\propto \sqrt{N} \ln(N)$ for comparison.

$\sqrt{N} \ln(N)$, which is consistent with the increasing ε values. For the σ SA-ES with $\tau = 1/\sqrt{2N}$ in the Rastrigin landscape, the population size is underestimated in accordance with the negative values of ε . The growth rate is slightly above $\sqrt{N} \ln(N)$ for $P_s = 0.5$ as predicted from the decreasing values of ε . The growth rate for $P_s = 0.99$ differs slightly from the predicted behavior. For both the CSA-ES with $c = 1/\sqrt{N}$ and the σ SA-ES with $\tau = 1/\sqrt{8N}$, the values of ε are positive and approximately constant with N . The growth rate is approximately $\sqrt{N} \ln(N)$. The predicted values from the equations are an upper bound for sufficiently large N , except for the Bohachevsky landscape for $P_s = 0.99$. In this case, the scaling behavior differs. For Bohachevsky and $P_s = 0.99$, the experimental values display fluctuations, whereas the values for $P_s = 0.5$ follow a rather smooth course. This discrepancy can be explained by the fact that the increase of P_s within the vicinity of $P_s \approx 0.99$ becomes weaker, see Fig. 8, and that the data was acquired from a small number of experiments.

In the case of the Ackley function, only an upper bound for the population size can be given. This means that the success probability is expected to approach 1 when μ is larger as

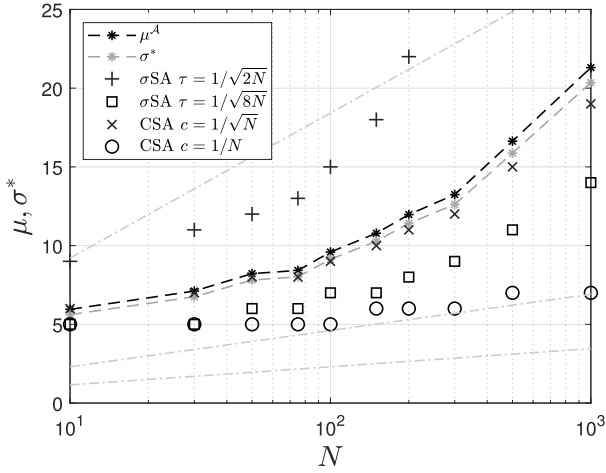


Fig. 14. Population size μ_{99} to reach $P_s = 99\%$ in the Ackley landscape with $\vartheta = 0.5$ and initialization values $R^{(0)} = 10\sqrt{N}$ and $\sigma^{*(0)} = 3$. Grey stars with dashed lines show σ^* experimentally determined for the CSA-ES with $c = 1/\sqrt{N}$ and $\mu = \mu_{99}$. Black stars show μ^A (64) using this σ^* . Gray dashed-dotted lines show functions $\propto \ln(N)$ for comparison.

predicted by (64). The dashed line with black stars in Fig. 14 shows (64), where σ^* was determined experimentally for the CSA-ES with $c = 1/\sqrt{N}$ and $\mu = \mu_{99}$. The values of σ^* , μ^A , and μ_{99} are almost identical and show the same scaling behavior. When using other strategy types or varying ϑ or the Ackley parameters, differences may be larger. In all cases, the scaling behaviors of σ^* and μ^A is approximately but slightly larger than μ_{99} .

Except for the σ SA-ES with $\tau = 1/\sqrt{2N}$, the experimental values increase slower than linearly. For the σ SA-ES with $\tau = 1/\sqrt{2N}$, the values of μ are significantly larger and increase faster than linearly. As demonstrated in Section V-E, the choice of suitable initial values depends on the strategy-specific parameters. By choosing different initial values, this anomalous behavior can be avoided and a sublinear scaling behavior can also be achieved.

VII. CONCLUSION

Each of the multimodal functions considered in this work can be divided into a global part, where the global optimizer is the only minimum, and a residual part with several local minima. The latter part can be interpreted as frozen noise. By applying the noisy sphere theory, a model has been developed for the calculation of the success probability P_s to reach the global optimizer of the multimodal function. P_s depends on the population size parameters μ and λ . For a noisy sphere with a fixed noise strength, an ES with a fixed population size cannot reach the optimizer arbitrarily close, but its parents fluctuate about the optimizer with an expected distance R_{st} . These fluctuations can be interpreted as a global search behavior. If R_{st} is sufficiently small the parents hit the global attractor region. In this region, the frozen noise wanes and the ES converges to the global optimizer. This convergence model is independent of specific ES strategy parameters and therefore applicable to different types of ES, such as σ SA or CSA-ES. The model was applied to the Rastrigin and Bohachevsky

function, with the result that the parental population size needed to achieve a satisfactory convergence probability scales approximately with $\mathcal{O}(\sqrt{N} \ln(N))$.

Even though the Bohachevsky test function formula appears more complex than that of Rastrigin, the analysis in this article revealed strong similarities w.r.t. the population sizing behavior. Therefore, from viewpoint of testbed design for empirical performance evaluations it suffices to include only one of these test functions. This should be taken into account when proposing new benchmark competitions. Similar considerations do also hold for the Ackley function.

As for the Ackley function, the frozen noise model does not allow for the determination of the success probability as there is no analytical description for the global attractor domain. Therefore, a different approach was used.

An upper bound for the population size μ was derived that mainly depends on σ^* . Since σ^* increases with N , the upper bound for the population size also depends on N , but only weakly. This is in contrast to the Rastrigin or Bohachevsky case. However, convergence depends strongly on the initialization. If the start population is located too far from the global optimizer the ES will only converge locally or diverge. This dependency from the initial conditions must be taken more into account when running empirical performance comparisons of evolutionary algorithms. The standard initializations found in literature appear to ensure convergence with a reasonable success probability for small population sizes (see Fig. 11). From this perspective, the standard Ackley function is not a hard test function. However, it could be made harder by decreasing the C_2 parameter in (3). This would elevate the local $\tilde{\mu}$ maximum in Fig. 12.

In addition to the multimodal functions considered, there is Griewank's function that can also be divided into a global part and a residual part. However, the Griewank function is an exception in the class of these multimodal functions because the success probability increases with N . The main reason for this counterintuitive behavior is the fast vanishing of the noise variance with N . That is, the problem gets easier with increasing N . From the viewpoint of analysis, a second peculiarity seems to complicate the application of the model additionally. The noise cannot be well approximated by a Gaussian distribution. Both peculiarities make the analysis in the interesting low-dimensional N cases appear questionable: The model developed in this article is based on $N \rightarrow \infty$ asymptotics. Thus, yielding only an approximation for small N that does not perform well for N of the order of 10 or below. Therefore, the Griewank case needs perhaps another approach and remains as a topic for another paper.

The investigations presented in this article do also contribute to a better understanding of how ES perform search and locate the optimizers in highly multimodal fitness landscapes. Contrary to popular belief, the ES does not follow a gradient-like path. The global optimizer is approached from all sides. That is, if correctly designed the ES performs a global search. However, this also means that there must be a global structure of the fitness landscape that guides the search of the ES in a highly noisy environment.

The analysis of ES in highly multimodal fitness landscapes is still in its very first beginning. One may question whether the results presented and methods developed are useful for real-world problems. There is no definite answer, yet. This question defines a research program on its own. Do the multimodal test functions to be found in the BBOB and CEC competitions really reflect real-world problems? Furthermore, progress in a theory-driven analysis will be always gradual. A possible direction could regard multimodal test functions where the global part cannot be well approximated by a sphere model. Also the performance of the ES on dynamically moving optimizer problems could define a research direction. And on the algorithm engineering side, the question of controlling the population size efficiently deserves attention.

ACKNOWLEDGMENT

For open access purposes, the author has applied a CCBY public copyright license to any author accepted manuscript version arising from this submission. The authors gratefully thank Amir Omeradzic for providing valuable feedback.

REFERENCES

- [1] N. Hansen and S. Kern, "Evaluating the CMA evolution strategy on multimodal test functions," in *Parallel Problem Solving From Nature*, vol. 8. Heidelberg, Germany: Springer, 2004, pp. 282–291.
- [2] N. Hansen, S. Müller, and P. Koumoutsakos, "Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES)," *Evol. Comput.*, vol. 11, no. 1, pp. 1–18, 2003.
- [3] L. Schönenberger and H.-G. Beyer, "On a population sizing model for evolution strategies optimizing the highly multimodal Rastrigin function," in *Proc. GECCO*, New York, NY, USA, 2023, pp. 848–855.
- [4] S. Meyer-Nieberg, "Self-adaptation in evolution strategies," Ph.D. dissertation, Dept. Comput. Sci., Univ. Dortmund, Dortmund, Germany, 2007.
- [5] N. Hansen, "Verallgemeinerte individuelle Schrittweitenregelung in der evolutionsstrategie," Ph.D. dissertation, Eine Untersuchung zur entstochastisierten, Tech. Univ. Berlin, Berlin, Germany, 1998.
- [6] D. Arnold, *Noisy Optimization With Evolution Strategies*. Dordrecht, The Netherlands: Kluwer Acad., 2002.
- [7] A. Omeradzic and H.-G. Beyer, "Progress rate analysis of evolution strategies on the Rastrigin function: First results," in *Parallel Problem Solving From Nature XXVII*, H. Aguirre et al., Eds., Berlin, Germany: Springer, 2022.
- [8] A. Omeradzic and H.-G. Beyer, "Progress analysis of a multi-recombinative evolution strategy on the highly multimodal Rastrigin function," *Theor. Comput. Sci.*, vol. 978, Nov. 2023, Art. no. 114179.
- [9] A. Omeradzic and H.-G. Beyer, "Convergence properties of the $(\mu/\mu_I, \lambda)$ -ES on the Rastrigin function," in *Proc. FOGA XVII*. ACM, 2023, pp. 117–128.
- [10] H.-G. Beyer, D. Arnold, and S. Meyer-Nieberg, "A new approach for predicting the final outcome of evolution strategy optimization under noise," *Genet. Program. Evol. Mach.*, vol. 6, no. 1, pp. 7–24, 2005.
- [11] H.-G. Beyer, "On the 'explorative power' of ES/EP-like algorithms," in *Proc. 7th Annu. Conf. Evol. Program. Evol. Program. VII*, 1998, pp. 323–334, doi: [10.1007/BFB0040785](https://doi.org/10.1007/BFB0040785).
- [12] D. Arnold and H.-G. Beyer, "Performance analysis of evolution strategies with multi-recombination in high-dimensional R^N -search spaces disturbed by noise," *Theor. Comput. Sci.*, vol. 289, pp. 629–647, Oct. 2002.
- [13] H.-G. Beyer, *The Theory of Evolution Strategies*. Heidelberg, Germany: Springer, 2001, doi: [10.1007/978-3-662-04378-3](https://doi.org/10.1007/978-3-662-04378-3).
- [14] P. J. Angeline, "Evolutionary optimization versus particle swarm optimization: Philosophy and performance differences," in *Proc. Int. Conf. Evol. Program.*, 1998, pp. 601–610.
- [15] G. Zhu and S. Kwong, "Gbest-guided artificial bee colony algorithm for numerical function optimization," *Appl. Math. Comput.*, vol. 217, no. 7, pp. 3166–3173, 2010.
- [16] B. Niu and H. Wang, "Bacterial colony optimization," in *Discrete Dynamics in Nature and Society*. Hoboken, NJ, USA: Wiley, 2012.



Lisa Schönenberger received the diploma degree in mathematics from the University of Vienna, Vienna, Austria, in 2010, and the bachelor's degree in computer science from FernUniversität Hagen, Hagen, Germany, in 2020.

Since 2021, she has been employed as a Research Assistant in the field of evolutionary algorithms, Vorarlberg University of Applied Sciences, Dornbirn, Austria. Her current research interests include evolutionary algorithms, multimodal optimization, and restart strategies.



Hans-Georg Beyer received the diploma degree in theoretical electrical engineering from Ilmenau Technical University, Ilmenau, Germany, in 1982, the Ph.D. degree in physics from Bauhaus-University Weimar, Weimar, Germany, in 1989, and the habilitation degree in computer science from the University of Dortmund, Dortmund, Germany, in 1997.

Since 2004, he has been a Professor with the Vorarlberg University of Applied Sciences, Dornbirn, Austria. He authored the book *The Theory of Evolution Strategies* (Heidelberg: Springer-Verlag, 2001) and authored/co-authored numerous papers in that field. He is best known for his theoretical analyses and the design of evolution strategies based on the stochastic dynamical systems approach.

Dr. Beyer was the Editor-in-Chief of the MIT Press Journal *Evolutionary Computation* and served as an Associate Editor for the IEEE TRANSACTIONS ON EVOLUTIONARY COMPUTATION from 1997 to 2021.