

UNIVERSITY OF DORTMUND

REIHE COMPUTATIONAL INTELLIGENCE

COLLABORATIVE RESEARCH CENTER 531

Design and Management of Complex Technical Processes
and Systems by means of Computational Intelligence Methods

Runtime Analysis of the (1+1) ES Minimizing
Simple Quadratic Forms Using the 1/5-Rule

Jens Jägersküpper

No. CI-186/04

Technical Report

ISSN 1433-3325

November 2004

Secretary of the SFB 531 · University of Dortmund · Dept. of Computer Science/XI
44221 Dortmund · Germany

This work is a product of the Collaborative Research Center 531, "Computational Intelligence," at the University of Dortmund and was printed with financial support of the Deutsche Forschungsgemeinschaft.

Runtime Analysis of the (1+1) ES Minimizing Simple Quadratic Forms Using the 1/5-Rule

Jens Jägersküpper*

FB Informatik, Lehrstuhl 2, Universität Dortmund, 44221 Dortmund, Germany
jj@Ls2.cs.uni-dortmund.de

Abstract. We consider the (1+1) Evolution Strategy, a simple evolutionary algorithm for continuous optimization problems, using so-called Gaussian mutations and the 1/5-rule for the adaptation of the mutation strength. Here, the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ to be minimized is given by a quadratic form $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$, where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a positive definite diagonal matrix and \mathbf{x} denotes the current search point. This is a natural extension of the well-known SPHERE-function ($\mathbf{Q}=\mathbf{I}$). Thus, very simple unconstrained quadratic programs are investigated, and the question is addressed how \mathbf{Q} effects the runtime. Therefore, quadratic forms

$$f(\mathbf{x}) = \xi \cdot (x_1^2 + \dots + x_{n/2}^2) + x_{n/2+1}^2 + \dots + x_n^2$$

with $\xi = \omega(1)$, i.e. $1/\xi \rightarrow 0$ as $n \rightarrow \infty$, and $\xi = \text{poly}(n)$ are exemplarily investigated. It is shown that the optimization very quickly approaches a steady state in which the expected runtime (defined as the number of f -evaluations) to halve the approximation error is $\Theta(\xi \cdot n)$. Though $\xi \cdot n = \text{poly}(n)$, this result actually shows that the evolving search point indeed creeps along the gentlest descent of the ellipsoidal fitness landscape.

1 Introduction

Finding—or at least approximating—an optimum of a given function $f: S \rightarrow \mathbb{R}$ is one of the fundamental problems—in theory as well as in practice. Methods for solving continuous optimization problems, e.g. $S = \mathbb{R}^n$, are usually classified into first-order, second-order, and zeroth-order methods depending on whether they utilize the gradient (the first derivative) of the objective function, the gradient and the Hessian (the second derivative), or none of both.¹

A zeroth-order method is also called *derivative-free* or *direct search method*. Newton’s method is the example for a second-order method; first-order methods can be (sub)classified into Quasi-Newton, steepest decent, and conjugate gradient methods. Classical zeroth-order methods try to approximate the gradient in order to plug this estimate into a first-order method. Finally, amongst the “modern” zeroth-order methods, evolutionary algorithms (EAs) come into play. EAs for continuous optimization, however, are usually subsumed under the term *evolution(ary)*

* supported by the German Research Foundation (DFG) as part of the research center “Computational Intelligence” (SFB 531)

¹ Note that here “continuous” relates to the search space rather than to f , and that, unlike in math programming, throughout this paper “ n ” denotes the number of dimensions of the search space and *not* the number of optimization steps; “ d ” generally denotes a distance in the n -dimensional search space.

strategies (ESs). Obviously, in general we cannot expect zeroth-order methods to out-perform first-order methods or even second-order methods.

However, when information about the gradient is not available, for instance if f relates to a property of some workpiece and is given by simulations or even by real-world experiments, first-order (and also second-order) methods just cannot be applied. As the approximation of the gradient usually involves $\Omega(n)$ f -evaluations, a single optimization step of a classical zeroth order-method is computationally intensive, especially if f is given implicitly by simulations. In practical optimization, especially in mechanical engineering, this is often the case, and particularly in this field EAs are enjoyed with growing popularity. However, the enthusiasm in practical EAs has led to an unclear variety of very sophisticated and problem-specific EAs. Unfortunately, from a theoretical point of view, the development of such EAs is solely driven by practical success and the aspect of a theoretical analysis is left aside. In other words,—concerning EAs—theory has not kept up with practice, and thus, we should not try to analyze the algorithmic runtime of the most sophisticated EA en vogue, but concentrate on very basic, or call them “simple”, EAs in order to build a sound and solid basis for EA-theory.

For discrete search spaces, essentially $\{0, 1\}^n$, such a theory has been started successfully in the mid 1990s (Mühlenbein (1992), cf. Wegener (2001) and Droste et al. (2002)). There first results for non-artificial but well-known problems have been obtained recently (namely for the maximum matching problem by Giel and Wegener (2003), for sorting and the shortest-path problem by Scharnow et al. (2002), and for the minimum-spanning tree by Neumann and Wegener (2004)).

The situation for continuous evolutionary optimization is different. Here, the vast majority of the results are based on empiricism, i. e., experiments are performed and their outcomes are interpreted, which leads to a theory in the sense of physics rather than computer science. Also convergence properties of EAs have been studied to a considerable extent (e. g. Rudolph (1997), Greenwood and Zhu (2001), Bienvenue and Francois (2003)). A lot of results have been obtained by analyzing a simplifying model of the stochastic process induced by the EA, for instance by letting the number of dimensions approach infinity. Unfortunately, such results rely on experimental validation as a justification for the simplifications/inaccuracies introduced by the modeling. In particular Beyer has obtained numerous results that focus on local performance measures (*progress rate, fitness gain*; cf. Beyer (2001b)), i. e., the effect of a single mutation (or, more generally, of a single transition from one generation to the next) is investigated. Best-case assumptions concerning the mutation adaptation in this single step then provide estimates of the maximum gain a single step may yield. However, when one aims at analyzing the (1+1) ES as an algorithm, rather than a model of the stochastic process induced, a different, more algorithmic approach is needed. In 2003 a first theoretical analysis of the expected runtime, given by the number of function evaluations, of the (1+1) ES using the 1/5-rule was presented (Jägersküpfer, 2003). The function/fitness landscape considered therein is the well-know SPHERE-function, given by $\text{SPHERE}(\mathbf{x}) := \sum_{i=1}^n x_i^2 = \mathbf{x}^\top \mathbf{I} \mathbf{x}$, and the multi-step behavior that the (1+1) ES bears when using the 1/5-rule for the adaptation of the mutation strength is rigorously analyzed. As mentioned in the abstract, the present paper will extend this result to a broader class of functions. One may guess that an ellipsoidal landscape is similar to the ridge-function scenario (especially to the parabolic ridge). Beyer (2001a) focuses on local measures for this

fitness landscape. However, since ridge functions are unbounded, i. e. there is no optimum, and there is no need for adaptation, from an algorithmic point of view—when one is interested in adaptation mechanisms and how they work—ellipsoidal fitness landscapes are more challenging.

Finally note that, regarding the approximation error, for unconstrained optimization it is generally not clear how the runtime can be measured (solely) with respect to the absolute error of the approximation. In contrast to discrete and finite problems, the initial error is generally not bounded, and hence, the question how many steps it takes to get into the ε -ball around an optimum does not make sense without specifying the starting conditions. Hence, we must consider the runtime with respect to the relative improvement of the approximation. Given that the optimization process becomes steady-state, considering the number of steps/ f -evaluations to halve the approximation error is a natural choice. For the SPHERE-function, Jägersküpfer (2003) gives a proof that the 1/5-rule makes the (1+1) ES perform $\Theta(n)$ steps to halve the distance from the optimum and, in addition, that this is asymptotically the best possible w. r. t. isotropically distributed mutation vectors, i. e., for any adaptation of isotropic mutations, the expected number of f -evaluations is $\Omega(n)$ (moreover, for any constant $\varepsilon > 0$, $O(n^{1-\varepsilon})$ f -evaluations suffice only with an exponentially small probability).

The Algorithm

We will concentrate on the (1+1) evolution strategy ((1+1) ES), which dates back to the mid 1960s (cf. Rechenberg (1973) and Schwefel (1995)). This simple EA uses solely mutation due to a single-individual population, where here “individual” is just a synonym for “search point”. Let $\mathbf{c} \in \mathbb{R}^n$ denote the current individual. Given a starting point, i. e. an initialization of \mathbf{c} , the (1+1) ES performs the following evolution loop:

1. Choose a random mutation vector $\mathbf{m} \in \mathbb{R}^n$, where the distribution of \mathbf{m} may depend on the course of the optimization process.
2. Generate the mutant $\mathbf{c}' \in \mathbb{R}^n$ by $\mathbf{c}' := \mathbf{c} + \mathbf{m}$.
3. IF $f(\mathbf{c}') \leq f(\mathbf{c})$ THEN \mathbf{c}' becomes the current individual ($\mathbf{c} := \mathbf{c}'$) ELSE \mathbf{c}' is discarded (\mathbf{c} unchanged).
4. IF the stopping criterion is met THEN output \mathbf{c} ELSE goto 1.

Since a worse mutant (with respect to the function to be minimized) is always discarded, the (1+1) ES is a randomized hill climber, and the selection rule is called *elitist selection*. Fortunately, for the type of results we are after we need not define a reasonable stopping criterion. How the mutation vectors are generated must be specified, though. Originally, the mutation vector $\mathbf{m} \in \mathbb{R}^n$ is generated by firstly generating a vector $\widetilde{\mathbf{m}} \in \mathbb{R}^n$ each component of which is independently standard normal distributed; subsequently, this vector is scaled by the multiplication with a scalar $s \in \mathbb{R}_{>0}$, i. e. $\mathbf{m} = s \cdot \widetilde{\mathbf{m}}$. This type of mutation is called *Gaussian mutation*. In practice, Gaussian mutations are the most common type of mutations (for the search space \mathbb{R}^n) and, therefore, will be considered here. The crucial property of a Gaussian mutation is that \mathbf{m} is isotropically distributed, i. e., $\mathbf{m}/|\mathbf{m}|$ is uniformly

distributed upon the unit hypersphere and the length of the mutation, namely the random variable $|\mathbf{m}|$, is independent of the direction $\mathbf{m}/|\mathbf{m}|$.²

The question that naturally arises is how the scaling factor s is to be chosen. Obviously, the smaller the approximation error, i. e., the closer \mathbf{c} is to an optimum, the shorter \mathbf{m} needs to be for a further improvement of the approximation to be possible. Unfortunately, the algorithm does not know about the current approximation error, but can utilize only the knowledge obtained by f -evaluations. Based on experiments and rough calculations for two function scenarios (namely SPHERE and a corridor function), Rechenberg proposed the *1/5-(success-)rule*. The idea behind this adaptation mechanism is that in a step of the (1+1) ES the mutant should be accepted with probability 1/5. Hereinafter, a mutation that results in $f(\mathbf{c}') \leq f(\mathbf{c})$ is called *successful*, and hence, when talking about a mutation, *success probability* denotes the probability that the mutant $\mathbf{c}' = \mathbf{c} + \mathbf{m}$ is at least as fit as \mathbf{c} . Obviously, when elitist selection is used, the success probability of a step equals the probability that the mutation is accepted in this step. If every step was successful with probability 1/5, we would observe that on the average one fifth of the mutations are successful. Thus, the 1/5-rule works as follows: the optimization process is observed for n steps without changing s ; if more than one fifth of the steps in this observation phase have been successful, s is doubled, otherwise, s is halved.³

The Function Scenario

In this section we will have a closer look to the fitness landscape under consideration and preview isotropic mutations in this scenario. Note that “fitness” is used as a synonym for “function value”. Furthermore, since the optimum function value is 0, the current approximation error is defined as $f(\mathbf{c})$, the fitness of the current individual. As mentioned in the abstract, we exemplarily consider the following class of quadratic forms ($n \in 2\mathbb{N}$):

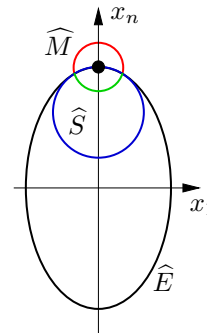
$$f_n(\mathbf{x}) := \xi \cdot (x_1^2 + \dots + x_{n/2}^2) + x_{n/2+1}^2 + \dots + x_n^2$$

Hence, $f_n(\mathbf{x}) = \xi \cdot \text{SPHERE}_{n/2}(\mathbf{y}) + \text{SPHERE}_{n/2}(\mathbf{z})$ where $\mathbf{y} := (x_1, \dots, x_{n/2})$ and $\mathbf{z} := (x_{n/2+1}, \dots, x_n)$. Thus, the aim is to minimize the sum of two separate sphere functions, in $S_1 = \mathbb{R}^{n/2}$ resp. $S_2 = \mathbb{R}^{n/2}$, having weight ξ resp. 1, i. e., $f(\mathbf{x}) = \xi \cdot |\mathbf{y}|^2 + |\mathbf{z}|^2$, where $|\cdot|$ denotes the length of a vector in Euclidean space (Euclidean norm). Recall that the mutation vector \mathbf{m} equals $s \cdot \widetilde{\mathbf{m}}$. As each component of $\widetilde{\mathbf{m}}$ is independently standard normal distributed, $\mathbf{m}_1 := (m_1, \dots, m_{n/2})$ and $\mathbf{m}_2 := (m_{n/2+1}, \dots, m_n)$ are two independent $(n/2)$ -dimensional Gaussian mutations based on the common scaling factor s . Obviously, \mathbf{m}_1 only affects \mathbf{y} , whereas \mathbf{m}_2 only affects \mathbf{z} , and thus, the fitness of the mutant equals $\xi \cdot |\mathbf{y} + \mathbf{m}_1|^2 + |\mathbf{z} + \mathbf{m}_2|^2$.

² The state of the art in mutation adaptation seems to be the *covariance matrix adaptation (CMA)* (Hansen and Ostermeier, 1996) where $\mathbf{m} = s \cdot \mathbf{B} \cdot \widetilde{\mathbf{m}}$ with a matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ which is also adapted. Unlike $\mathbf{B} = t \cdot \mathbf{I}$ for some scalar t , the mutation vector is not isotropically distributed—by intention, of course.

³ Various implementations of the 1/5-rule can be found in the literature, yet in fact, one result of (Jägersküpfer, 2003) is that the order of the runtime is indeed not affected as long as the observation lasts $\Theta(n)$ steps and the scaling factor s is multiplied by a constant greater than 1 resp. by a positive constant smaller than 1.

Let $d_1 := |\mathbf{y}|$ and $d_2 := |\mathbf{z}|$ denote the distance from the origin/optimum in S_1 resp. S_2 . Since Gaussian mutations as well as SPHERE are invariant with respect to rotations of the coordinate system, we may rotate S_1 and S_2 such that we are located at $(d_1, 0, \dots, 0) \in S_1$ resp. $(0, \dots, 0, d_2) \in S_2$. In other words, we may assume w. l. o. g. that the current search point is located at $(d_1, 0, \dots, 0, d_2) \in \mathbb{R}^n$, i. e., that it lies in the x_1 - x_n -plane. In fact, we have just described a projection $\hat{\cdot}: \mathbb{R}^n \rightarrow \mathbb{R}^2$. Note that due to the properties of f and Gaussian mutations this projection only conceals irrelevant information, i. e., all information relevant to the analysis is preserved. Thus, we can concentrate on the 2D-projection as depicted in the figure. For some arguments, however, it is crucial to keep in mind that this projection is based on the fact that the current search point, and also its mutant, can be assumed to lie in the x_1 - x_n -plane w. l. o. g. (obviously, for the mutant to lie in this plane, S_1 and S_2 must almost surely (a. s.) be re-rotated).



In the next section some of the results presented in (Jägersküpper, 2003), which will be used here, will be shortly restated. In Section 3 the crucial properties of a single mutation in the considered fitness landscape are discussed, and in the subsequent section we will have a closer look to the adaptation, i. e., the multi-step behavior of the (1+1) ES will be analyzed for the considered function class/fitness landscape. We end with some concluding remarks in Section 5.

2 Preliminaries

In this section some notions and notations are introduced. Furthermore, the results obtained for the SPHERE-scenario in (Jägersküpper, 2003) that we will use are resumed; for more details cf. Jägersküpper (2002).

Definition 1. A probability $p(n)$ is **exponentially small** in n if for a constant $\varepsilon > 0$, $p(n) = \exp(-\Omega(n^\varepsilon))$. An event $A(n)$ happens **with overwhelming probability (w. o. p.)** with respect to n if $\mathbb{P}\{\neg A(n)\}$ is exponentially small in n .

A statement $Z(n)$ holds **for n large enough** if $(\exists n_0 \in \mathbb{N})(\forall n \geq n_0) Z(n)$.

Let $\mathbf{c} \in \mathbb{R}^n - \{\mathbf{0}\}$ denote a search point and \mathbf{m} a Gauss mutation. Note that $\text{SPHERE}(\mathbf{c}) = |\mathbf{c}|^2$ (recall that $|\mathbf{c}|$ is the L^2 -norm (Euclidian length) of \mathbf{c}). The analysis of the (1+1) ES for SPHERE has shown that for n large enough

$$\mathbb{P}\{|\mathbf{c} + \mathbf{m}| \leq |\mathbf{c}| \mid |\mathbf{m}| = \ell\} \geq \varepsilon \text{ for a constant } \varepsilon \in (0, \frac{1}{2}) \iff \ell = O(|\mathbf{c}|/\sqrt{n}),$$

i. e., the mutant obtained by an isotropic mutation of \mathbf{c} is closer to a predefined point, here the origin, with probability $\Omega(1)$ iff the length of the mutation vector is at most an $O(1/\sqrt{n})$ -fraction of the distance between \mathbf{c} and this point. On the other hand,

$$\mathbb{P}\{|\mathbf{c} + \mathbf{m}| \leq |\mathbf{c}| \mid |\mathbf{m}| = \ell\} \leq \varepsilon \text{ for a constant } \varepsilon \in (0, \frac{1}{2}) \iff \ell = \Omega(|\mathbf{c}|/\sqrt{n}),$$

in other words, the mutant obtained by an isotropic mutation of \mathbf{c} is closer to a predefined point, here again the origin, with a constant probability strictly smaller than $1/2$ iff the length of the mutation vector is at least an $\Omega(1/\sqrt{n})$ -fraction of the

distance between \mathbf{c} and this point. (The actual constant ε respectively correlates with the (multiplicative) constant in the O -notation resp. in the Ω -notation.)

The expected length of a Gauss mutation \mathbf{m} equals $s \cdot \mathbb{E}[|\widetilde{\mathbf{m}}|] = s \cdot \Theta(\sqrt{n})$ since $|\widetilde{\mathbf{m}}|$ is χ -distributed (with n degrees of freedom). Let $\bar{\ell} := \mathbb{E}[|\mathbf{m}|]$. Moreover, $\mathbb{P}\{||\mathbf{m}| - \bar{\ell}| \geq \delta \cdot \bar{\ell}\} \leq 1/(\delta^2(2n-1))$ for $\alpha > 0$, in other words, there is only small deviation in the length of a mutation; e. g., with probability $1 - O(1/n)$ the mutation vector's actual length differs from its expected length by no more than $\pm 1\%$.

Concerning the mutation adaptation by the 1/5-rule for SPHERE, we know that there exists a constant $p_h \in (0, 1/5)$ such that if the success probability of the mutation in the first step of an observation phase is smaller than p_h , then w. o. p. less than 1/5 of the steps in this phase are successful so that the scaling factor is halved. Analogously, $p_d \in (1/5, 1/2)$ exists such that if the first step of a phase is successful with probability at least p_d , then w. o. p. more than 1/5 of the steps in this phase are successful so that s is doubled. This can be used to show that the 1/5-rule in fact ensures that each step is successful with a probability in $[a, b] \subset (0, 1/2)$ for two constants a, b .

Let $\Delta = |\mathbf{c}| - |\mathbf{c}'|$ denote the spatial gain towards the origin, the optimum of SPHERE, in a step. For SPHERE, a mutation is accepted (by elitist selection) iff $\Delta \geq 0$. Consequently, negative gains are zeroed out. Thus, the expected spatial gain of a step is $\mathbb{E}[\Delta \cdot \mathbb{1}_{\{\Delta \geq 0\}}]$ and we know that $\mathbb{E}[\Delta \cdot \mathbb{1}_{\{\Delta \geq 0\}} \mid \bar{\ell} = \Theta(|\mathbf{c}|/\sqrt{n})] = \Theta(|\mathbf{c}|/n)$, i. e., if the scaling factor causes a mutation to be successful with a constant probability in $(0, 1/2)$ —for instance 1/5—then the distance from the optimum is expected to decrease by an $\Theta(1/n)$ -fraction. Furthermore, in this situation the distance decreases by an $\Omega(1/n)$ -fraction with probability $\Omega(1)$.

3 Gain in a Single Step

In this section we now take a closer look at the properties of a Gaussian mutation in the ellipsoidal fitness landscape under consideration. Since $\xi = \omega(1)$, $\xi > 1$ for n large enough, and therefore, we assume $\xi > 1$ in the following. Furthermore, “ f ” will also be used as an abbreviation of the fitness of the current individual and “ f' ” stands for the fitness of the mutant.

Recall that $f = \xi d_1^2 + d_2^2$ (for the current search point) and $f' = \xi d_1'^2 + d_2'^2$ (for its mutant), where $d_1' := |\mathbf{y} + \mathbf{m}_1|$ and $d_2' := |\mathbf{z} + \mathbf{m}_1|$. The crucial point to the analysis is the answer to the question how d_1 , d_2 and the scaling factor s —and with it $|\mathbf{m}|$ —relate when the success probability of a step, i. e. the probability that the mutant is accepted, is about 1/5. In other words, how does the length of the mutation vector depend on d_1 and d_2 , and how do d_1 and d_2 relate. Since $\nabla \widehat{f}(d_1, d_2) = (\xi 2 d_1, 2 d_2)^\top$, for a search point satisfying $d_1/d_2 = 1/\xi$ an infinitesimal change of d_1 has the same effect on f as an infinitesimal change of d_2 . Though the length of a mutation is not infinitesimal, this may be taken as an indicator that the ratio d_1/d_2 will approach a steady state when using isotropic mutations, and indeed, it turns out that the process approaches a steady state where $d_1/d_2 = \Theta(1/\xi)$. In this section, we will see that near the gentlest descent in our ellipsoidal fitness landscape, namely for $d_1/d_2 = O(1/\xi)$, a mutation succeeds with a constant probability greater than 0 but smaller than 1/2 iff the scaling factor corresponds to $\mathbb{E}[|\mathbf{m}|] = \Theta(\sqrt{f/n}/\xi)$. Furthermore, asymptotically tight bounds on the expected fitness gain of a single step, i. e. $\mathbb{E}[f - \min\{f, f'\}]$, in such a situ-

ation will be obtained. Therefore, we will show that a mutation of a search point \mathbf{c} for which $d_1/d_2 = O(1/\xi)$ with a mutation satisfying $\mathbb{E}[|\mathbf{m}|] = \Theta(\sqrt{f/n}/\xi)$ in the ellipsoidal fitness landscape is “similar” to the mutation (with the same scaling factor) of a search point \mathbf{x} in the SPHERE scenario with $\text{SPHERE}(\mathbf{x}) = \Theta(f/\xi^2)$.

We start our analysis with $d_1 = 0$ and $d_2 = \phi$ so that $f = \phi^2$. Consequently, \mathbf{c} is located at a point with gentlest descent w.r.t. all points having fitness ϕ^2 , and hence, the curvature of the 2D-curve given by the projection of $f = \phi^2$, i.e. the projection \widehat{E} of the n -ellipsoid $E := \{\mathbf{x} \mid f(\mathbf{x}) = f(\mathbf{c})\} \subset \mathbb{R}^n$, is maximum at $\widehat{\mathbf{c}}$. By a simple application of differential geometry (Appendix A), we get that the curvature equals ξ/ϕ (for $d_1 = 0$ and $d_2 = \phi$). Consequently, the radius of the osculating circle (\widehat{S} in the figure) equals ϕ/ξ . As this circle \widehat{S} actually lies in the x_1 - x_n -plane, it is the equator of an n -sphere S with radius ϕ/ξ (the center of which lies on the x_n -axis, just like the current search point \mathbf{c}). Note that this sphere lies completely inside E such that $S \cap E = \{\mathbf{c}\}$. Thus, the probability that a mutation hits inside S is a lower bound on the probability that $f' \leq f$, i.e.,

$$\begin{aligned} \mathbb{P}\{f' \leq f\} &= \mathbb{P}\{\mathbf{c} + \mathbf{m} \text{ lies inside } E\} \\ &\geq \mathbb{P}\{\mathbf{c} + \mathbf{m} \text{ lies inside } S\} \\ &= \mathbb{P}\left\{|\mathbf{x} + \mathbf{m}| \leq |\mathbf{x}| \text{ for some } \mathbf{x} \text{ with } |\mathbf{x}| = \text{radius of } \widehat{S} = \phi/\xi\right\} \\ &= \mathbb{P}\{\text{SPHERE}(\mathbf{x} + \mathbf{m}) \leq \text{SPHERE}(\mathbf{x}) \mid \text{SPHERE}(\mathbf{x}) = (\phi/\xi)^2\}. \end{aligned}$$

In fact, our argumentation yields that the preceding (in)equalities hold for any fixed length ℓ of the mutation vector \mathbf{m} , i.e., if the probabilities are conditioned on the event $\{|\mathbf{m}| = \ell\}$, respectively. Since ℓ is arbitrary here and the radius of S is independent of ℓ , they remain valid when this condition is dropped.

For an upper bound on the probability that a mutation hits inside E , consider a mutation (vector) having length $\ell < 2\phi$ (since for $\ell \geq 2\phi$, E lies inside M). Let $M = \{\mathbf{x} \mid |\mathbf{c} - \mathbf{x}| = \ell\} \subset \mathbb{R}^n$ denote the mutation sphere consisting of all potential mutants. Then \widehat{M} is a circle (cf. the figure above) with radius ℓ centered at $\widehat{\mathbf{c}}$. (Note that, though $\mathbf{c}' = \mathbf{c} + \mathbf{m}$, given $|\mathbf{m}| = \ell$, is uniformly distributed upon M , $\widehat{\mathbf{c}'}$ is *not* uniformly distributed upon \widehat{M}). Now consider the curvature at a point in $\widehat{E} \cap \widehat{M} = \{\mathbf{z}_1, \mathbf{z}_2\}$ (there are exactly two points of intersection since $0 < \ell < 2\phi$). Simple differential geometry shows that the curvature at \mathbf{z}_i is $\kappa_\ell = \Theta(\xi/\phi)$ if $\ell = O(\phi/\xi)$ (cf. Appendix A). As the curvature at any point of \widehat{E} that lies inside \widehat{M} is greater than κ_ℓ (since $\xi > 1$), $\widehat{\mathbf{c}}$ as well as \mathbf{z}_i lie inside the osculating circle at \mathbf{z}_{3-i} which has radius $r_\ell := 1/\kappa_\ell = \Theta(\phi/\xi)$ if $\ell = O(\phi/\xi)$. Thus, there is also a circle with radius r_ℓ passing through $\widehat{\mathbf{c}}$ such that \mathbf{z}_1 and \mathbf{z}_2 lie inside this circle. Therefore, the circle passing through \mathbf{z}_1 , \mathbf{z}_2 , and $\widehat{\mathbf{c}}$ has a radius smaller than r_ℓ , and again, this circle actually lies in the x_1 - x_n -plane of the search space and is the image of the n -sphere having this circle as an equator. Hence,

$$\begin{aligned} &\mathbb{P}\{f' \leq f \mid |\mathbf{m}| = \ell\} \\ &\leq \mathbb{P}\{\text{SPHERE}(\mathbf{x} + \mathbf{m}) \leq \text{SPHERE}(\mathbf{x}) \mid \text{SPHERE}(\mathbf{x}) = (\alpha \phi/\xi)^2, |\mathbf{m}| = \ell\} \end{aligned}$$

where $\alpha = \Theta(1)$ if $\ell = O(\phi/\xi)$. (Besides, $r_\ell \searrow \phi/\xi$ as $\ell \searrow 0$.)

Recall that we assumed $\widehat{\mathbf{c}} = (0, \phi) \in \mathbb{R}^2$, i.e. $d_1 = 0$ and $d_2 = \phi$, in the above argumentation. The estimates we have made for the bounds on the probability of a mutation hitting inside the n -ellipsoid E , however, remain valid as long as

$d_1/d_2 = O(1/\xi)$: Since ξ/ϕ is the maximum curvature of \widehat{E} , there is always a circle \widehat{S} with radius ϕ/ξ lying inside \widehat{E} such that $\widehat{S} \cap \widehat{E} = \{\widehat{\mathbf{c}}\}$, and since \widehat{S} is in fact an equator of an n -sphere S , S lies completely inside E such that $S \cap E = \{\mathbf{c}\}$. For the upper bound, we must merely consider the \mathbf{z}_i at which the curvature is smaller, and indeed, it turns out that as long as $d_1/d_2 = O(1/\xi)$ and $\ell = O(\phi/\xi)$, κ_ℓ is still $\Theta(\xi/\phi)$ (cf. Appendix A).

Hence, when $f(\mathbf{c}) = \phi^2$ such that \mathbf{c} satisfies $d_1/d_2 = O(1/\xi)$, we are in a situation resembling (w. r. t. the success probability of a mutation) the minimization of SPHERE at a point having distance $\Theta(\phi/\xi)$ from the optimum/origin. Concerning the 1/5-rule, we know (cf. Section 2) that

$$(\exists \text{ constant } \varepsilon > 0) \varepsilon \leq \mathbb{P}\{f' \leq f\} \leq 1/2 - \varepsilon \iff \ell = \Theta((\phi/\xi)/\sqrt{n})$$

where ε correlates with the two multiplicative constants within the Θ -notation.

Thus, we are now going to investigate the gain of a step when $f = \phi^2$ and $\ell = \Theta((\phi/\xi)/\sqrt{n})$. As we have seen above, there exists an n -sphere S with radius $r = \phi/\xi$ lying completely in E such that $S \cap E = \{\mathbf{c}\}$. Again owing to the results for SPHERE, we know that a mutation having length $\ell = r/\sqrt{n}$ hits with probability $\Omega(1)$ a hyperspherical cap $C \subset M$ containing all points of M that are at least $\Omega(r/n)$ closer to the center of S than \mathbf{c} . Consequently, with probability $\Omega(1)$ the mutant lies inside E such that its distance from E is $\Theta(r/n)$, i. e. $\Theta((\phi/\xi)/n)$. If we pessimistically assume that this spatial gain were realized along the gentlest descent of f , i. e. $d_1 = 0$ and $d'_1 = 0$ so that $d'_2 = d_2 - \Theta((\phi/\xi)/n)$, we obtain that with probability $\Omega(1)$

$$\begin{aligned} f' &\leq (\phi - \Theta((\phi/\xi)/n))^2 \\ &= \phi^2 - 2\alpha\phi^2/(\xi n) + \alpha^2\phi^2/(\xi n)^2 \text{ for some } \alpha = \Theta(1) \\ &= \phi^2 - \underbrace{\alpha(2 - \alpha/(\xi n))}_{\Theta(1)} \phi^2/(\xi n) \\ &= \phi^2 - \Theta(1) \phi^2/(\xi n) \\ &= f - \Theta(f/(\xi n)). \end{aligned}$$

Let $\mathbf{c}'' := \arg \min\{f(\mathbf{c}), f(\mathbf{c}')\}$ denote the search point that gets selected by elitist selection. Since mutants with a worse fitness are rejected, i. e. $f'' \leq f$, this implies for the expected fitness gain of a step

$$\mathbb{E}\left[f'' \mid |\mathbf{m}| = \Theta(\sqrt{f/n}/\xi)\right] = f - \Omega(f/(\xi n)).$$

Due to the pessimistic assumptions, this lower bound on the fitness gain just derived is valid only for $\ell = \Theta(\sqrt{f/n}/\xi)$, yet it holds independently of the ratio d_1/d_2 . A spatial gain of $\Theta(f/(\xi n))$ could result in a much larger fitness gain, though. If $d_1/d_2 = O(1/\xi)$, however, the fitness gain is also $O(f/(\xi n))$ as we will see. Therefore, let $d_1 = \alpha \cdot \phi/\xi$ with $\alpha = O(1)$ and still $f = \xi \cdot d_1^2 + d_2^2 = \phi^2$. Owing to the argumentation for the upper bound on the success probability of a step, we know that there is an n -sphere S with radius $r = O(\phi/\xi)$ such that $\mathbf{c} \in S$ and $I := M \cap E \subset S$, where I is the boundary of the hyperspherical cap $C \subset M$ lying inside E . Owing to the results for SPHERE, we know that $\mathbb{E}[\text{dist}(\mathbf{c}', I) \cdot \mathbf{1}_{\{\mathbf{c}' \in C\}}] = O(r/n)$ even if $|\mathbf{m}|$ is optimal, i. e., even if the length of the mutation vector were magically chosen such that the expected distance of the selected search point, \mathbf{c}'' , from the center of S

is minimized. In other words, we know that if a mutation hits inside E , its expected distance from E is $O(r/n) = O((\phi/\xi)/n)$ anyway. Thus, if we optimistically assume that the spatial gain were realized completely in S_1 , i. e. completely on the heavier weighted SPHERE, (so that $d'_2 = d_2$, implying $d''_2 = d_2$), we obtain

$$\begin{aligned} \mathbb{E}[\xi d_1''^2 + d_2''^2 \mid d_1/d_2 = O(1/\xi)] &\geq \xi (d_1 - O((\phi/\xi)/n))^2 + d_2^2 \\ &= \xi (\alpha\phi/\xi - O((\phi/\xi)/n))^2 + d_2^2 \\ &\geq \xi ((\alpha\phi/\xi)^2 - 2\alpha(\phi/\xi) \cdot O((\phi/\xi)/n)) + d_2^2 \\ &= \xi d_1^2 - O(\phi^2/(\xi n)) + d_2^2 \end{aligned}$$

and hence,

$$\mathbb{E}[f'' \mid d_1/d_2 = O(1/\xi)] = \phi^2 - O(\phi^2/(\xi n)) = f - O(f/(\xi n)).$$

This upper bound on the expected fitness gain of a step holds only for $d_1/d_2 = O(1/\xi)$, yet independently of (the distribution of) $|\mathbf{m}|$, which is converse to the lower bound. However, altogether we have proved the following:

Lemma 1. *Consider a step of the (1+1) ES. If $d_1/d_2 = O(1/\xi)$ in this step, then $\varepsilon \leq \mathbb{P}\{f' \leq f\} \leq 1/2 - \varepsilon$ for a constant $\varepsilon > 0$ iff $|\mathbf{m}| = \Theta(\sqrt{f/n}/\xi)$.*

If $d_1/d_2 = O(1/\xi)$ and $|\mathbf{m}| = \Theta(\sqrt{f/n}/\xi)$ in this step, then $\mathbb{E}[f - f''] = \Theta(f/(\xi n))$, and furthermore, $f - f'' = \Omega(f/(\xi n))$ with probability $\Omega(1)$.

4 Multi-Step Behavior

The results just obtained show that if $d_1/d_2 = O(1/\xi)$ during a phase of n steps (an observation phase of the 1/5-rule) and $\mathbb{E}[|\mathbf{m}|] = \Theta(\sqrt{f/n}/\xi)$, i. e. $\varepsilon \leq \mathbb{P}\{f' \leq f\} \leq 1/2 - \varepsilon$ for a constant $\varepsilon > 0$, at the beginning of this phase, then we expect $\Theta(n)$ steps each of which reduces the fitness by $\Theta(f/(\xi n))$. By Chernoff bounds, there are $\Omega(n)$ such steps w. o. p., and thus, the fitness, and with it the approximation error, is reduced w. o. p. by an $\Theta(1/\xi)$ -fraction in this phase. Consequently, after $\Theta(\xi)$ consecutive phases, the fitness is halved w. o. p. if during these phases $d_1/d_2 = O(1/\xi)$. Since, up to now, the argumentation completely bases on the results for SPHERE, even the argumentation on the 1/5-rule can be adapted, which directly yields the following result (cf. Theorem 2 in (Jägersküpper, 2003) or Theorem 3 in Jägersküpper (2002)):

Theorem 1. *If $d_1/d_2 = O(1/\xi)$ in the complete optimization process and the scaling factor is initialized such that in the first step $\mathbb{E}[|\mathbf{m}|] = \Theta(\sqrt{f/n}/\xi)$, then the expected number of steps/ f -evaluations to reduce the initial approximation error to a 2^{-t} -fraction, $t = \text{poly}(n)$, is $\Theta(t \cdot \xi \cdot n)$.*

Obviously, the assumption “ $d_1/d_2 = O(1/\xi)$ in the complete optimization process” lacks justification and is, therefore, objectionable. It must be replaced by a much more weaker assumption on the starting conditions only. Thus, the crucial point in the analysis is the question why should $d_1/d_2 = O(1/\xi)$. This question will be tackled in the following. Therefore, let $\Delta_1 := d_1 - d'_1$ and $\Delta_2 := d_2 - d'_2$ denote the spatial gain of the mutant towards the origin in S_1 resp. S_2 . Then

$$f' = \xi (d_1 - \Delta_1)^2 + (d_2 - \Delta_2)^2 = \xi d_1^2 - \xi 2d_1\Delta_1 + \xi \Delta_1^2 + d_2^2 - 2d_2\Delta_2 + \Delta_2^2,$$

and hence,

$$f' \leq f \iff f' - f \leq 0 \iff -\xi 2d_1 \Delta_1 + \xi \Delta_1^2 - 2d_2 \Delta_2 + \Delta_2^2 \leq 0.$$

Let α be defined by $\alpha/\xi = d_1/d_2$. Then the latter inequality is equivalent to

$$\begin{aligned} -2\alpha d_2 \Delta_1 + \xi \Delta_1^2 - 2d_2 \Delta_2 + \Delta_2^2 \leq 0 &\iff -\alpha \Delta_1 + \frac{\xi \Delta_1^2}{2d_2} \leq \Delta_2 - \frac{\Delta_2^2}{2d_2} \\ \text{(using } d_2 = \xi \cdot d_1/\alpha) &\iff -\alpha \Delta_1 \left(1 - \frac{\Delta_1}{2d_1}\right) \leq \Delta_2 \left(1 - \frac{\Delta_2}{2d_2}\right) \end{aligned}$$

Thus, when using elitist selection, the mutant is accepted iff the last inequality holds. Note that whenever a mutation satisfying $-\alpha \Delta_1 > \Delta_2$ is accepted, then

$$1 - \frac{\Delta_1}{2d_1} < 1 - \frac{\Delta_2}{2d_2} \iff \frac{\Delta_1}{d_1} > \frac{\Delta_2}{d_2} \iff \Delta_1 > \frac{d_1}{d_2} \Delta_2 \iff \Delta_1 > \frac{\alpha}{\xi} \Delta_2,$$

implying that $\Delta_1 > 0$ and $\Delta_2 < 0$, and consequently, such a step surely results in $d_1''/d_2'' < d_1/d_2$, i. e. $\alpha'' < \alpha$. Hence, in the following we may concentrate on the accepted mutations for which $-\alpha \Delta_1 \leq \Delta_2$.

So, let us assume for a moment that the mutant replaces/becomes the current individual iff $-\alpha \Delta_1 \leq \Delta_2$. Let $i \in \{1, 2\}$. As Δ_{3-i} is random, $\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]$ is a random variable taking the value $\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq x\}}]$ whenever Δ_2 happens to take the value x . We are interested in $\mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]] = \mathbb{E}[d_i - d_i'']$, the expected reduction of the distance from the optimum in S_i in a step, and we expect $d_1''/d_2'' \leq d_1/d_2$, i. e. $\alpha'' \leq \alpha$, iff

$$\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]] \geq \frac{\alpha}{\xi} \cdot \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]].$$

In order to prove that this inequality holds for $\alpha = O(1)$, we aim at a lower bound on $\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]]$ and an upper bound on $\mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]]$ in the following. Note that

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]] &= \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i < 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} < 0\}}] \\ &\quad + \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i < 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} \geq 0\}}] \\ &\quad + \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} < 0\}}] \\ &\quad + \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} \geq 0\}}] \end{aligned}$$

and that $\mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i < 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} < 0\}}] = 0$ since the three indicator inequalities describe the empty set. Since $\Delta_1, \Delta_2 \geq 0$ implies $-\alpha \Delta_1 \leq \Delta_2$,

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_i \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} \geq 0\}}] &= \mathbb{E}[\mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{\Delta_i \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_{3-i} \geq 0\}}] \\ &= \mathbb{E}[\Delta_i \cdot \mathbb{1}_{\{\Delta_i \geq 0\}}] \cdot \mathbb{P}\{\Delta_{3-i} \geq 0\}. \end{aligned}$$

As we need a lower bound on $\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}}]]$, we may pessimistically assume that $\Delta_1 = -x/\alpha$ whenever Δ_2 happens to equal x . By this assumption,

$$\begin{aligned} &\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_1 < 0\}}] \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \\ &\geq -\mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha \Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_1 < 0\}}]/\alpha, \end{aligned}$$

$$\begin{aligned} & \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \\ & \geq -\mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] / \alpha. \end{aligned}$$

All in all, we have

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] & \geq \mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\} \\ & \quad - \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_1 < 0\}}] / \alpha \\ & \quad - \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] / \alpha, \end{aligned}$$

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] & = \mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{P}\{\Delta_1 \geq 0\} \\ & \quad + \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_1 < 0\}}] \\ & \quad + \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}]. \end{aligned}$$

Recall that we want to show that for some $\alpha = O(1)$

$$\xi \cdot \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \geq \alpha \cdot \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]],$$

and note that $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\}$ and $\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{P}\{\Delta_1 \geq 0\}$ are of the same order when $\mathbb{P}\{\Delta_2 \geq 0\}$ and $\mathbb{P}\{\Delta_1 \geq 0\}$ are $\Omega(1)$, respectively. Consequently, since $\xi = \omega(1)$, for the above inequality to hold for n large enough, it would be sufficient that

$$\begin{aligned} & \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_1 < 0\}}] \\ & + \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \leq 0 \end{aligned} \quad (1)$$

because then we would have

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] & \geq \mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\} \quad \text{and} \\ \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] & \leq \mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{P}\{\Delta_1 \geq 0\}. \end{aligned}$$

Concerning the expected spatial gain in S_2 , however, we are going to use the trivial upper bound

$$\mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \leq \mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}],$$

and thus, we concentrate on a lower bound on the expected spatial gain in S_1 in the following. Therefore, we prove next that inequality (1) holds for $\alpha = O(1)$ at least if the actual length of \mathbf{m}_2 differs by no more than a constant fraction from $\bar{\ell}_1$, the expected length of \mathbf{m}_1 .

Lemma 2. *If $\mathbb{P}\{\Delta_1 \geq 0\} = \Omega(1)$ and $|\mathbf{m}_2| = \Theta(\bar{\ell}_1)$, there exists a constant α^* such that for n large enough inequality (1) holds for all $\alpha \geq \alpha^*$.*

The proof can be found in Appendix B. Note that $\bar{\ell}_1 = \bar{\ell}_2$ in our scenario. We know (cf. Section 2) that

$$\mathbb{P}\left\{ \left| |\mathbf{m}_2| - \bar{\ell}_2 \right| \geq (\sqrt{3} - 1) \cdot \bar{\ell}_2 \right\} \leq \left((\sqrt{3} - 1)^2 \cdot 2 \cdot (n - 1) \right)^{-1} < (n - 1)^{-1},$$

and thus, the condition “ $|\mathbf{m}_2| = \Theta(\bar{\ell}_1)$ ” is violated only with probability $O(1/n)$. Whether or not this condition is met, obviously $\Delta_1 \geq -|\mathbf{m}_1|$, and consequently, $\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \geq -\bar{\ell}_1$. Applying this rough/trivial bound only in case of $||\mathbf{m}_2| - \bar{\ell}_1| > (\sqrt{3} - 1) \cdot \bar{\ell}_1$ and $\Delta_1 < 0 \leq \Delta_2 \vee \Delta_1 \geq 0 > \Delta_2$, we can extend the preceding lemma: if $\mathbb{P}\{\Delta_1 \geq 0\} = \Omega(1)$ then for $\alpha \geq \alpha^*$

$$\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \geq \mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\} - \frac{\bar{\ell}_1}{n-1}.$$

Next we will see that this additive error term vanishes in situations that arise due to the 1/5-rule.

Lemma 3. *If $\mathbb{P}\{\Delta_1 \geq 0\}$, $1/2 - \mathbb{P}\{\Delta_1 \geq 0\}$, $\mathbb{P}\{\Delta_2 \geq 0\}$ are $\Omega(1)$, respectively, there exists a constant α^* such that for $\alpha \geq \alpha^*$ and n large enough*

$$\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]] \geq \mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\} / 2.$$

Proof. Recall that $f' \leq f \wedge -\alpha\Delta_1 > \Delta_2$ implies $\Delta_1 > 0 > \Delta_2$. Consequently, all (Δ_1, Δ_2) -tuples zeroed out by $\mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}$, but kept by $\mathbb{1}_{\{f' \leq f\}}$ are in $\mathbb{R}_{>0} \times \mathbb{R}_{<0}$. Analogously, $f' > f \wedge -\alpha\Delta_1 \leq \Delta_2$ implies $\Delta_1 < 0 < \Delta_2$ so that all (Δ_1, Δ_2) -tuples kept by $\mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}$, but zeroed out by $\mathbb{1}_{\{f' \leq f\}}$ are in $\mathbb{R}_{<0} \times \mathbb{R}_{>0}$. Hence,

$$\begin{aligned} \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]] &\geq \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \\ (\text{and } \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{f' \leq f\}}]]) &\leq \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]]). \end{aligned}$$

As $1/2 - \mathbb{P}\{\Delta_1 \geq 0\}, \mathbb{P}\{\Delta_1 \geq 0\} = \Omega(1)$ implies $\bar{\ell}_1 = \Theta(d_1/\sqrt{n})$ and, as a consequence, $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] = \Theta(d_1/n)$ (cf. Section 2), the error term $\bar{\ell}_1/(n-1)$ is $\Theta(d_1/n^{1.5})$ whereas $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\} = \Theta(d_1/n)$. As the error term is by a $\Theta(1/\sqrt{n})$ -factor smaller, finally $1 - \Theta(1/\sqrt{n}) \geq 1/2$ for n large enough. \square

Recall: we expect $\alpha'' = \alpha$ iff $\xi \cdot \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]] = \alpha \cdot \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{f' \leq f\}}]]$ or, equivalently, iff $\mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]]/d_1 = \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{f' \leq f\}}]]/d_2$. Thus there exists a distinct α_0 such that there is no drift w. r. t. the ratio d_1/d_2 , and this is the steady state of the optimization: for $\alpha < \alpha_0$, α is more likely to increase than to decrease, and for $\alpha > \alpha_0$, α is more likely to decrease than to increase.

Since $\mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{f' \leq f\}}]] \leq \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}}]] \leq \mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}]$ and $\xi = \omega(1)$, we have $\xi \cdot \mathbb{P}\{\Delta_2 \geq 0\}/2 \geq \alpha^*$ for n large enough if $\mathbb{P}\{\Delta_2 \geq 0\} = \Omega(1)$, and hence, $\alpha_0 \leq \alpha^* = O(1)$ under the conditions of Lemma 3. Besides, the 1/5-rule just ensures these conditions as long as $d_1 = O(d_2)$. For the same reasons, there exists $\alpha_\downarrow > \alpha_0$ such that $\xi \cdot \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]] \geq 2 \cdot \alpha \cdot \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{f' \leq f\}}]]$ (for n large enough) and $\alpha_\downarrow = O(1)$ again under the conditions of Lemma 3. Thus, when $\alpha \geq \alpha_\downarrow$ there is a drift towards smaller α ; more formally:

Lemma 4. *Let the scaling factor s be fixed. If $\mathbb{P}\{\Delta_1 \geq 0\}$ and $1/2 - \mathbb{P}\{\Delta_1 \geq 0\}$ as well as $\mathbb{P}\{\Delta_2 \geq 0\}$ are $\Omega(1)$, respectively, there exists a constant α_\downarrow such that for n large enough, if in the i^{th} step $\alpha^{[i]} \geq \alpha_\downarrow$ (yet $\alpha^{[i]} = O(\xi)$), then w. o. p. after at most $n^{0.3}$ steps the search is located at a point for which $\alpha < \alpha^{[i]}$, and furthermore, w. o. p. $\alpha \leq \alpha^{[i]} + O(\alpha^{[i]}/n^{0.6})$ in all intermediate steps.*

The proof can be found in Appendix C. Since the 1/5-rule keeps the scaling factor unchanged for n steps, we can virtually partition each such observation phase in

$n/n^{0.3} = n^{0.7}$ sub-phases to each of which this lemma applies. Incorporating these new insights into the argumentation for the 1/5-rule known from the analysis of SPHERE finally enables us to replace the objectionable condition “ $d_1/d_2 = O(1/\xi)$ in the complete optimization process” in Theorem 1 by “ $d_1/d_2 = O(1/\xi)$ for the initial search point” —yielding the main result on the steady state performance of the the (1+1) ES on the considered quadratic forms:

Theorem 2. *If $d_1/d_2 = O(1/\xi)$ for the initial search point and the scaling factor is initialized such that in the first step $\mathbb{E}[|\mathbf{m}|] = \Theta(\sqrt{f/n}/\xi)$, then the expected number of steps/ f -evaluations to reduce the initial approximation error/function value to a 2^{-t} -fraction, $t = \text{poly}(n)$, is $\Theta(t \cdot \xi \cdot n)$.*

Now, one might ask what happens if the optimization starts at a point for which d_1 is not $O(d_2/\xi)$. A closer look at the argumentation in the proof of the preceding lemma reveals that the same argumentation results in the proof of the existence of another constant $\alpha_{\Downarrow} > \alpha_{\downarrow}$ such that the drift towards smaller α is that strong when $\alpha \geq \alpha_{\Downarrow}$ that w. o. p. α drops by a constant fraction within at most n steps:

Lemma 5. *Let the scaling factor s be fixed. If $\mathbb{P}\{\Delta_1 \geq 0\}$ and $1/2 - \mathbb{P}\{\Delta_1 \geq 0\}$ as well as $\mathbb{P}\{\Delta_2 \geq 0\}$ are $\Omega(1)$, respectively, then there exists a constant α_{\Downarrow} such that for n large enough: if in the i^{th} step $\alpha^{[i]} \geq \alpha_{\Downarrow}$ (yet $\alpha^{[i]} = O(\xi)$), then w. o. p. after at most n steps the search is located at a point with $\alpha \leq \alpha^{[i]} - \Omega(\alpha^{[i]})$.*

See Appendix D for the proof. Finally, this lemma shows that α drops very quickly if the lemma’s conditions are met. Again utilizing the results for SPHERE, it is rather simple to check that these conditions are met when d_1 is $O(d_2)$ (and $\Omega(d_2/\xi)$, of course). If d_1 is not $O(d_2)$, for instance if we start at a point of steepest descent, i. e. $d_2 = 0$ so that $f = \xi d_1^2$, then a simple argumentation using rough bounds on Δ_1 and Δ_2 yields that—as expected— d_1/d_2 drops even faster than in situations covered by the preceding lemma since the (expected) spatial gain in S_1 (on the heavier weighed sphere) is negative whereas the one in S_2 is positive.

5 Conclusion

Based on the results on how the (1+1) ES minimizes the well-known SPHERE-function, we have extended these results to a broader class of functions consisting of certain positive definite quadratic forms. The main insight of the results presented is that Gaussian mutations adapted by the 1/5-rule result in the optimization process to become steady-state very close to the gentlest descent of the ellipsoidal fitness landscape. However, more insight into how EAs for continuous optimization work is gained, contributing to building an algorithmic EA-theory for continuous search spaces.

A The Curvature is $\Omega(\xi/\phi)$ if $d_1 = O(d_2/\xi)$

We consider the ellipse given by $\xi \cdot d_1^2 + d_2^2 = \phi^2$. Thus, $d_2 = \sqrt{\phi^2 - \xi \cdot d_1^2}$,

$$\frac{dd_2}{dd_1} = \frac{-\xi \cdot d_1}{\sqrt{\phi^2 - \xi \cdot d_1^2}} \quad \text{and} \quad = \frac{dd_2}{d^2 d_1} = \frac{-\xi^2 \cdot d_1^2}{(\phi^2 - \xi \cdot d_1^2)^{3/2}} + \frac{-\xi}{\sqrt{\phi^2 - \xi \cdot d_1^2}}.$$

As the curvature (of a plane curve) equals

$$\frac{\frac{dd_2}{d^2d_1}}{\left(1 + \left(\frac{dd_2}{d^2d_1}\right)^2\right)^{3/2}} = \frac{\phi^2\xi}{(\phi^2 + (\xi^2 - \xi) \cdot d_1^2)^{3/2}},$$

for $d_1 = \alpha \cdot \phi/\xi$ the curvature equals

$$\frac{\xi}{\phi \cdot (1 + (1 - 1/\xi) \cdot \alpha^2)^{3/2}}.$$

Finally note that $(1 + (1 - 1/\xi) \cdot \alpha^2)^{3/2} = O(1)$ for $\alpha = O(1)$, i. e. $d_1 = O(\phi/\xi)$. Furthermore, this shows that for $d_1 = 0$, i. e. $\alpha = 0$, the curvature equals ξ/ϕ .

B Proof of Lemma 2

Proof. Let us assume for a moment that the distribution of $|\mathbf{m}_2|$ were concentrated at a certain ℓ_2 , and let “ $D\{\cdot\}$ ” denote the density of an event. Then

$$\begin{aligned} & \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 \geq 0\}}] \cdot \mathbb{1}_{\{\Delta_1 < 0\}}] \\ &= \int_0^{\ell_2} x \cdot D\{\Delta_2 = x\} \cdot \mathbb{P}\{-x/\alpha \leq \Delta_1 < 0\} dx \quad \text{and} \\ & \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbb{1}_{\{-\alpha\Delta_1 \leq \Delta_2\}} \cdot \mathbb{1}_{\{\Delta_2 < 0\}}] \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \\ &= \int_{-\ell_2}^0 y \cdot D\{\Delta_2 = y\} \cdot \mathbb{P}\{\Delta_1 \geq -y/\alpha\} dy \\ &= \int_0^{\ell_2} -x \cdot D\{\Delta_2 = -x\} \cdot \mathbb{P}\{\Delta_1 \geq x/\alpha\} dx. \end{aligned}$$

We know from the analysis of SPHERE that for $x \in [0, \ell_2)$

$$D\{\Delta_2 = x\} < \frac{\Psi_n}{\ell_2} \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} < D\{\Delta_2 = -x\}$$

(with $\Psi_n := \pi^{-1/2} \cdot \Gamma(n/2) / \Gamma(n/2 - 1/2) = \Theta(\sqrt{n})$, where “ Γ ” denotes the Gamma function), and thus, the LHS of (1) is smaller than

$$\begin{aligned} & \int_0^{\ell_2} x \cdot \frac{\Psi_n}{\ell_2} \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \mathbb{P}\{-x/\alpha \leq \Delta_1 < 0\} dx \\ & - \int_0^{\ell_2} x \cdot \frac{\Psi_n}{\ell_2} \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \mathbb{P}\{\Delta_1 \geq x/\alpha\} dx \\ &= \int_0^{\ell_2} x \cdot \frac{\Psi_n}{\ell_2} \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot (\mathbb{P}\{-x/\alpha \leq \Delta_1 < 0\} - \mathbb{P}\{\Delta_1 \geq x/\alpha\}) dx. \end{aligned}$$

Let $\Phi: [0, \ell_2] \rightarrow [-1, 1]$ be defined by $\Phi(y) := \mathbb{P}\{-y \leq \Delta_1 < 0\} - \mathbb{P}\{\Delta_1 \geq y\}$. Hence,

$$\int_0^{\ell_2} x \cdot \frac{\Psi_n}{\ell_2} \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \Phi(x/\alpha) dx \leq 0$$

implies the inequality (1). Note that, obviously, $\mathbf{P}\{-0 \leq \Delta_1 < 0\} = 0$ and, by assumption, $\mathbf{P}\{\Delta_1 \geq 0\} = \Omega(1)$. Since, $\mathbf{P}\{\Delta_1 \geq y\}$ decreases monotonically, whereas $\mathbf{P}\{-y \leq \Delta_1 < 0\}$ increases monotonically when y grows, $\Phi(y)$ is monotone increasing for $0 \leq y \leq \min\{\ell_1, \ell_2\}$ and equals $\mathbf{P}\{\Delta_1 < 0\}$ for $y \geq \ell_1$. Furthermore, if ε denotes an arbitrary constant with $0 < \varepsilon < \mathbf{P}\{\Delta_1 \geq 0\}$, then $\mathbf{P}\{\Delta_1 \geq y\} = \varepsilon$ implies $y = \Theta(\bar{\ell}_1/\sqrt{n})$. Analogously, if $0 < \varepsilon < \mathbf{P}\{\Delta_1 < 0\}$, then $\mathbf{P}\{-y \leq \Delta_1 < 0\} = \varepsilon$ implies $y = \Theta(\bar{\ell}_1/\sqrt{n})$. Thus, there exists $\check{y} = \kappa \cdot \bar{\ell}_1/\sqrt{n-1}$ with $\kappa = \Theta(1)$ such that $\mathbf{P}\{\Delta_1 \geq \check{y}\} = \mathbf{P}\{-\check{y} \leq \Delta_1 < 0\}$, i. e., $\Phi(\check{y}) = 0$, and hence, the inequality to be shown reads

$$\begin{aligned} & -\frac{\Psi_n}{\ell_2} \int_0^{\alpha \cdot \check{y}} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \Phi(x/\alpha) dx \\ \geq & \frac{\Psi_n}{\ell_2} \int_{\alpha \cdot \check{y}}^{\ell_2} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \Phi(x/\alpha) dx. \end{aligned} \quad (2)$$

For the RHS we have, using $(1 - a/(n-1))^{(n-1)/2} \leq e^{-a/2}$ for $n-1 > a > 0$,

$$\begin{aligned} & \int_{\alpha \cdot \check{y}}^{\ell_2} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \Phi(x/\alpha) dx \\ & \leq \int_{\alpha \cdot \check{y}}^{\ell_2} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot 1 dx \\ & = \left[\frac{-\ell_2^2}{2} \cdot \frac{(1 - (x/\ell_2)^2)^{(n-1)/2}}{(n-1)/2} \right]_{\alpha \cdot \check{y}}^{\ell_2} \\ & = 0 - \left(\frac{-\ell_2^2}{n-1} \cdot (1 - (\alpha \cdot \check{y}/\ell_2)^2)^{(n-1)/2} \right) \\ & = \frac{\ell_2^2}{n-1} \cdot (1 - (\alpha \cdot \check{y}/\ell_2)^2)^{(n-1)/2} \\ & \leq \frac{\ell_2^2}{n-1} \cdot (1 - (\alpha \cdot \kappa \cdot \bar{\ell}_1/\ell_2)^2/(n-1))^{(n-1)/2} \\ \text{if } n-1 > \left(\alpha \cdot \kappa \cdot \frac{\bar{\ell}_1}{\ell_2} \right)^2 \text{ then } & \leq \frac{\ell_2^2}{n-1} \cdot e^{-(\alpha \cdot \kappa \cdot \bar{\ell}_1/\ell_2)^2/2}. \end{aligned}$$

For the LHS of (2) note that, by the same arguments, there exists $\check{y} = \tau \cdot \bar{\ell}_1/\sqrt{n-1}$ with $\tau = \Theta(1)$ such that $\mathbf{P}\{\Delta_1 \geq \check{y}\} = 2 \cdot \mathbf{P}\{-\check{y} \leq \Delta_1 < 0\}$, and thus, for $0 \leq y \leq \check{y}$ we have $\mathbf{P}\{\Delta_1 \geq y\} \geq 2 \cdot \mathbf{P}\{-y \leq \Delta_1 < 0\}$, i. e., $-\Phi(y) \geq p := \mathbf{P}\{\Delta_1 \geq \check{y}\}/2 = \Omega(1)$. Hence,

$$\begin{aligned} & -\int_0^{\alpha \cdot \check{y}} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot \Phi(x/\alpha) dx \\ \geq & \int_0^{\alpha \cdot \check{y}} x \cdot (1 - (x/\ell_2)^2)^{(n-3)/2} \cdot p dx \\ & = p \cdot \left[\frac{-\ell_2^2}{2} \cdot \frac{(1 - (x/\ell_2)^2)^{(n-1)/2}}{(n-1)/2} \right]_0^{\alpha \cdot \check{y}} \\ & = p \cdot \frac{-\ell_2^2}{n-1} \cdot \left((1 - (\alpha \cdot \check{y}/\ell_2)^2)^{(n-1)/2} - 1 \right) \\ & = p \cdot \frac{\ell_2^2}{n-1} \cdot \left(1 - \left(1 - \frac{(\alpha \cdot \tau \cdot \bar{\ell}_1/\ell_2)^2}{n-1} \right)^{(n-1)/2} \right) \\ \text{if } n-1 > \left(\alpha \cdot \tau \cdot \frac{\bar{\ell}_1}{\ell_2} \right)^2 \text{ then } & \geq p \cdot \frac{\ell_2^2}{n-1} \cdot \left(1 - e^{-(\alpha \cdot \tau \cdot \bar{\ell}_1/\ell_2)^2/2} \right) \end{aligned}$$

All in all, we have broken it down into the inequality

$$p \cdot \frac{\ell_2^2}{n-1} \cdot \left(1 - e^{-(\alpha \cdot \tau \cdot \bar{\ell}_1 / \ell_2)^2 / 2}\right) \geq \frac{\ell_2^2}{n-1} \cdot e^{-(\alpha \cdot \kappa \cdot \bar{\ell}_1 / \ell_2)^2 / 2}.$$

Since p , τ , and κ are $\Theta(1)$, it is finally obvious that $\alpha = O(1)$ can be chosen large enough for this inequality to hold for n large enough if $\bar{\ell}_1 / \ell_2 = \Theta(1)$, i. e. $\ell_2 = \Theta(\bar{\ell}_1)$. \square

C Proof of Lemma 4

Proof. We begin by proving the second claim. Let us assume that, starting with the i^{th} step, $\alpha \geq \alpha^{[i]}$ for $k \leq n^{0.3}$ steps. Recall that, due to elitist selection, the fitness is non-increasing. As $d_2 > d_2^{[i]}$ and $f \leq f^{[i]}$ implies $d_1 < d_1^{[i]}$, which again implies $\alpha / \xi = d_1 / d_2 < d_1^{[i]} / d_2^{[i]} = \alpha^{[i]} / \xi$, we have just proved that (surely) $d_2 \leq d_2^{[i]}$ in these k steps, respectively. Since, irrespective of the adaptation of the length of an isotropic mutation, in a step w. o. p. $\Delta_2 = O(d_2 / n^{0.9})$, in all $k \leq n^{0.3}$ steps w. o. p. $d_2 \geq d_2^{[i]} - k \cdot O(d_2^{[i]} / n^{0.9}) \geq d_2^{[i]} - O(d_2^{[i]} / n^{0.6})$, i. e., $d_2 = d_2^{[i]}(1 - \sigma)$ for some $\sigma = O(n^{-0.6})$, respectively. Concerning an upper bound on d_1 , we have

$$f = \xi d_1^2 + d_2^2 = \xi d_1^2 + \left(d_2^{[i]} - \sigma d_2^{[i]}\right)^2 \leq f^{[i]} = \xi d_1^{[i]2} + d_2^{[i]2},$$

and hence

$$\begin{aligned} \xi d_1^2 &\leq \xi d_1^{[i]2} + (2\sigma - \sigma^2) d_2^{[i]2} \\ \Leftrightarrow d_1^2 &\leq d_1^{[i]2} + (2\sigma - \sigma^2) \frac{d_2^{[i]2}}{\xi} = d_1^{[i]2} + (2\sigma - \sigma^2) \frac{d_1^{[i]2}}{\alpha^{[i]}} \\ &= d_1^{[i]2} \left(1 + \frac{\sigma(2 - \sigma)}{\alpha^{[i]}}\right) \end{aligned}$$

Since $\sigma(2 - \sigma) / \alpha^{[i]}$ is $O(\sigma)$, i. e. $O(n^{-0.6})$, we finally get that in all k steps

$$\frac{\alpha}{\xi} = \frac{d_1}{d_2} \leq \frac{d_1^{[i]}}{d_2^{[i]}} \cdot \frac{\sqrt{1 + O(n^{-0.6})}}{1 - O(n^{-0.6})} = \frac{\alpha^{[i]}}{\xi} \cdot (1 + O(n^{-0.6})).$$

Now we are ready for the proof of the lemma's first claim. Therefore, assume that $\alpha \geq \alpha^{[i]} \geq \alpha_{\downarrow}$ for $n^{0.3} + 1$ steps. We are going to show that the probability of observing such a sequence of steps is exponentially small. Note that, since w. o. p. $d_2 \geq d_2^{[i]}(1 - \sigma)$ as we have seen, our assumption implies that also w. o. p. $d_1 \geq d_1^{[i]}(1 - \sigma)$, i. e., w. o. p. $d_1 = d_1^{[i]} - O(d_1^{[i]} / n^{0.6})$ in all $n^{0.3}$ steps. Let $X_j^{[k]}$, $j \in \{1, 2\}$, denote the RV $\Delta_j \cdot \mathbf{1}_{\{f' \leq f\}}$ in the $(i - 1 + k)^{\text{th}}$ step (so that $\mathbb{E}[X_j] = \mathbb{E}[\mathbb{E}[\Delta_j \cdot \mathbf{1}_{\{f' \leq f\}}]]$). Then, according to the arguments preceding the lemma, for $1 \leq k \leq n^{0.3}$, $\mathbb{E}[X_1^{[k]}] / d_1^{[k]} \geq 2 \cdot \mathbb{E}[X_2^{[k]}] / d_2^{[k]}$, i. e.,

$$\xi \cdot \mathbb{E}[X_1^{[k]}] \geq 2 \cdot \alpha^{[k]} \cdot \mathbb{E}[X_2^{[k]}] \geq 2 \cdot \alpha^{[i]} \cdot \mathbb{E}[X_2^{[k]}].$$

Let $S_j^{[k]} := X_j^{[1]} + \dots + X_j^{[k]}$ denote the total gain of k steps w. r. t. to d_j . By linearity of expectation, $\mathbb{E}[S_1^{[k]}/d_1^{[i]}] \geq 2 \cdot \mathbb{E}[S_2^{[k]}/d_2^{[i]}]$ for $1 \leq k \leq n^{0.3}$; however, the goal is to show that $\mathbb{P}\{S_1^{[k]}/d_1^{[i]} \leq S_2^{[k]}/d_2^{[i]} \text{ for } 1 \leq k \leq n^{0.3}\}$ is exponentially small.

Therefore, we will assume the worst case (w. r. t. to the analysis, i. e. the best case w. r. t. the chance of observing such a sequence) that $\mathbb{E}[X_1^{[k]}/d_1^{[i]}] = 2 \cdot \mathbb{E}[X_2^{[k]}/d_2^{[i]}]$ in each step. To see that this is in fact the worst case consider a search point \mathbf{x} for which $\alpha \geq \alpha^{[i]}$, i. e. $d_1/d_2 > d_1^{[i]}/d_2^{[i]}$, so that $\xi \cdot \mathbb{E}[X_1] > 2 \cdot \alpha \cdot \mathbb{E}[X_2]$. Now consider a search point $\tilde{\mathbf{x}}$ with $f(\tilde{\mathbf{x}}) = f(\mathbf{x})$ but $\tilde{\alpha} < \alpha$, i. e., $\tilde{d}_1 < d_1$ and $\tilde{d}_2 > d_2$. Owing to the results on SPHERE we know that, for an isotropic mutation of an arbitrary fixed length ℓ_j , for any fixed $g \in (-\ell_j, \ell_j)$, $\mathbb{P}\{\Delta_j \geq g\}$ strictly increases with d_j (when $d_j > \ell_j$). Consequently, (independently of the distribution of $|\mathbf{m}|$) $\tilde{\Delta}_1$ is stochastically dominated by Δ_1 , whereas $\tilde{\Delta}_2$ stochastically dominates Δ_2 . This implies that X_1 dominates \tilde{X}_1 , whereas X_2 is dominated by \tilde{X}_2 (in particular, we have $\mathbb{E}[X_1] < \mathbb{E}[\tilde{X}_1]$ and $\mathbb{E}[X_2] > \mathbb{E}[\tilde{X}_2]$).

As we have just seen, we may pessimistically assume that in each step the search is located at a point for which $\xi \cdot \mathbb{E}[X_1] = 2 \cdot \alpha \cdot \mathbb{E}[X_2]$. Hence, $\mathbb{E}[S_1^{[k]}/d_1^{[i]}] = 2 \cdot \mathbb{E}[S_2^{[k]}/d_2^{[i]}]$. Let $S_j := S_j^{[n^{0.3}]}$. Since $1.2/0.8 = 1.5 < 2$, it is sufficient to show that w. o. p. $S_1 \geq 0.8 \cdot \mathbb{E}[S_1]$ and w. o. p. $S_2 \leq 1.2 \cdot \mathbb{E}[S_2]$. The Hoeffding bounds (1963) (cf. Section 2.6.2 of (Hofri, 1987)) state that, for $X_j^{[k]} \in [a_j, b_j]$ and $t_j > 0$,

$$\begin{aligned} \mathbb{P}\{S_1 - \mathbb{E}[S_1] \leq -n^{0.3} \cdot t_1\} &\leq \exp\left(\frac{-2 \cdot n^{0.3} \cdot t_1^2}{(b_1 - a_1)^2}\right) \quad \text{and} \\ \mathbb{P}\{S_2 - \mathbb{E}[S_2] \geq n^{0.3} \cdot t_2\} &\leq \exp\left(\frac{-2 \cdot n^{0.3} \cdot t_2^2}{(b_2 - a_2)^2}\right). \end{aligned}$$

For $t_j = 0.2 \cdot \mathbb{E}[S_j]/n^{0.3}$, both exponents equal

$$-0.08 \cdot n^{-0.3} \cdot \mathbb{E}[S_j]^2 / (b_j - a_j)^2 = -\Omega(n^{-0.3}) \cdot \left(\frac{\mathbb{E}[S_j]}{b_j - a_j}\right)^2,$$

respectively. Therefore, our goal is to show that $\mathbb{E}[S_j]/(b_j - a_j) = \Omega(n^{0.2})$.

First we concentrate on $\mathbb{E}[S_1]$. Since S_1 is the sum of $n^{0.3}$ RVs $X_1^{[k]}$, it suffices to show that $\mathbb{E}[X_1^{[k]}/(b_1 - a_1)] = \Omega(n^{-0.1})$ for $1 \leq k \leq n^{0.3}$. In the following we assume that $d_1 = d_1^{[i]} \pm O(d_1^{[i]}/n^{0.6})$ and $d_2 \in [d_2^{[i]} - O(d_2^{[i]}/n^{0.6}), d_2^{[i]}]$ since we have seen (in the preceding proof of the second claim) that this happens w. o. p. Owing to the results for SPHERE, we know that $\mathbb{P}\{\Delta_j \geq 0\} = \Omega(1)$ implies that the scaling factor s is $O(d_j/n)$, which results in $\bar{\ell}_j = O(d_j/\sqrt{n})$, and that, under these conditions, w. o. p. $|\Delta_j| = O(\bar{\ell}_j/n^{0.4})$. Recall that $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{f' \leq f\}}]$ is at least $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\}/2$. Since $\mathbb{P}\{\Delta_2 \geq 0\} = \Omega(1)$ in i^{th} step and $d_2 \geq d_2^{[i]}(1 - O(n^{-0.6}))$ in all $n^{0.3}$ steps, in each of these steps $\mathbb{P}\{\Delta_2 \geq 0\} = \Omega(1)$. Hence, $\mathbb{E}[X_1] = \Omega(\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}])$ in each of the $n^{0.3}$ steps. Owing to the results for SPHERE, we know that (since $\bar{\ell}_1 = O(d_1/\sqrt{n})$ as we have seen) $\mathbb{E}[\Delta_1 \cdot \mathbb{1}_{\{\Delta_1 \geq 0\}}] = \Theta(\bar{\ell}_1/\sqrt{n})$ so that $\mathbb{E}[X_1] = \Omega(\bar{\ell}_1/\sqrt{n})$. Thus, $\mathbb{E}[S_1] = n^{0.3} \cdot \Omega(\bar{\ell}_1/\sqrt{n}) = \Omega(\bar{\ell}_1/n^{0.2})$ and $b_1 - a_1 = O(\bar{\ell}_1/n^{0.4})$, i. e., $\mathbb{E}[S_1]/(b_1 - a_1) = \Omega(n^{0.2})$.

Concerning a lower bound on $\mathbb{E}[S_2]$, recall that $\mathbb{E}[S_1]/d_1^{[i]} = 2 \cdot \mathbb{E}[S_2]/d_2^{[i]}$, i. e., $\mathbb{E}[S_2] = \mathbb{E}[S_1] \cdot d_2^{[i]}/(2 \cdot d_1^{[i]}) = \Omega(\bar{\ell}_1/n^{0.2}) \cdot \Omega(\xi/\alpha^{[i]})$. As $\bar{\ell}_1 = \bar{\ell}_2$ and (by assumption)

$\alpha^{[i]} = O(\xi)$, we have $\mathbb{E}[S_2] = \Omega(\bar{\ell}_2/n^{0.2})$. Since $b_2 - a_2 = O(\bar{\ell}_2/n^{0.4})$ (see above), $\mathbb{E}[S_2]/(b_2 - a_2) = \Omega(\bar{\ell}_2/n^{0.2})/O(\bar{\ell}_2/n^{0.4})$ is also $\Omega(n^{0.2})$.

All in all, our initial assumption that $\alpha \geq \alpha^{[i]} \geq \alpha_1$ for $n^{0.3} + 1$ steps implies that w. o. p. for the first $n^{0.3}$ steps $S_1/S_2 > \alpha^{[i]}/\xi$, i. e., that w. o. p. after at most $n^{0.3}$ steps α drops below $\alpha^{[i]}$ — showing that the sequence of steps we assumed to be observed happens only with an exponentially small probability. \square

D Proof of Lemma 5

Proof. By the same arguments used before, under the given assumptions there exists $\alpha' = O(1)$ such that for n large enough $\xi \cdot \mathbb{E}[\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{f' \leq f\}}]] \geq \mathbf{3} \cdot \alpha \cdot \mathbb{E}[\mathbb{E}[\Delta_2 \cdot \mathbf{1}_{\{f' \leq f\}}]]$. Let $\alpha_\Downarrow := 2 \cdot \alpha'$. Assume that $\alpha^{[i]} \geq \alpha_\Downarrow$ and $\alpha \geq \alpha_\Downarrow/2 = \alpha'$ for n steps (if $\alpha < \alpha_\Downarrow/2$ within these n steps, there is nothing to show). Following the same argumentation used in the proof of the preceding lemma (except for S_j now being the sum of n (instead of $n^{0.3}$) RVs), we get that w. o. p. $S_1/S_2 > 2 \cdot \alpha^{[i]}/\xi$, and hence, after these n steps w. o. p.

$$\begin{aligned} \frac{d_1}{d_2} &\leq \frac{d_1^{[i]} - S_1}{d_2^{[i]} - S_2} < \frac{d_1^{[i]} - S_1}{d_2^{[i]} - S_1 \cdot \xi/(2 \cdot \alpha^{[i]})} = \frac{d_1^{[i]} - S_1}{(d_1^{[i]} - S_1/2) \cdot \xi/\alpha^{[i]}} \\ &= \frac{d_2^{[i]} - S_1}{d_1^{[i]} - S_1/2} \cdot \frac{\alpha^{[i]}}{\xi} = \left(1 - \frac{S_1/2}{d_1^{[i]} - S_1/2}\right) \cdot \frac{d_1^{[i]}}{d_2^{[i]}} \leq \left(1 - \frac{S_1}{2 \cdot d_1^{[i]}}\right) \cdot \frac{d_1^{[i]}}{d_2^{[i]}}. \end{aligned}$$

Thus, we must merely show that $S_1 = \Omega(d_1^{[i]})$. Recall that S_1 is the sum of n RVs $X_1^{[k]}$ ($\Delta_1 \cdot \mathbf{1}_{\{f' \leq f\}}$ in the $(i - 1 + k)^{\text{th}}$ step, respectively). In the following we consider the i th step. Our argumentation just bases on the fact that $\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{f' \leq f\}}] \geq \mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{\Delta_1 \geq 0\}}] \cdot \mathbb{P}\{\Delta_2 \geq 0\}/2$, and since $\mathbb{P}\{\Delta_2 \geq 0\} = \Omega(1)$ by assumption, we have $\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{f' \leq f\}}] = \Omega(\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{\Delta_1 \geq 0\}}])$. Furthermore, since $\mathbb{P}\{\Delta_1 \geq 0\}$ as well as $1/2 - \mathbb{P}\{\Delta_1 \geq 0\}$ are $\Omega(1)$ by assumption, we know that $\bar{\ell}_1 = \Theta(d_1/\sqrt{n})$, which implies $\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{\Delta_1 \geq 0\}}] = \Theta(d_1/n)$. As a consequence, the assumptions ensure $\mathbb{E}[\Delta_1 \cdot \mathbf{1}_{\{f' \leq f\}}] = \Omega(d_1/n)$, and hence, $\mathbb{E}[S_1] = n \cdot \Omega(d_1/n) = \Omega(d_1)$. Applying Hoeffding's bound just as in the proof of the preceding lemma, we immediately get that w. o. p. $S_1 = \Omega(\mathbb{E}[S_1]) = \Omega(d_1^{[i]})$. \square

Bibliography

- Beyer, H.-G. [2001a]. On the performance of $(1, \lambda)$ -evolution strategies for the ridge function class. *IEEE Transactions on Evolutionary Computation* 5(3): 218–235.
- Beyer, H.-G. [2001b]. *The Theory of Evolution Strategies*. Springer.
- Bienvenue, A. and Francois, O. [2003]. Global convergence for evolution strategies in spherical problems: Some simple proofs and difficulties. *Theoretical Computer Science* 306: 269–289.
- Droste, S., Jansen, T., and Wegener, I. [2002]. On the analysis of the $(1+1)$ evolutionary algorithm. *Theoretical Computer Science* 276: 51–82.
- Giel, O. and Wegener, I. [2003]. Evolutionary algorithms and the maximum matching problem. *Proc. of the 20th Int'l Symposium on Theoretical Aspects of Computer Science (STACS)*, LNCS 2607. Springer, 415–426.
- Greenwood, G. W. and Zhu, Q. J. [2001]. Convergence in evolutionary programs with self-adaptation. *Evolutionary Computation* 9(2): 147–157.
- Hansen, N. and Ostermeier, A. [1996]. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. *Proc. of the IEEE Int'l Conference on Evolutionary Computation (ICEC)*. 312–317.
- Hoeffding, W. [1963]. Probability inequalities for sums of bounded random variables. *American Statistical Association Journal* 58(301): 13–30.
- Hofri, M. [1987]. *Probabilistic Analysis of Algorithms*. Springer.
- Jägersküpper, J. [2002]. Analysis of a simple evolutionary algorithm for the minimization in euclidian spaces. Technical Report CI-140/02, Univ. Dortmund, SFB 531. [http://sfbc1.uni-dortmund.de](http://sfbc1.uni-dortmund.de/Publications/Tech-Reports)→Publications→Tech-Reports.
- Jägersküpper, J. [2003]. Analysis of a simple evolutionary algorithm for minimization in Euclidian spaces. *Proc. of the 30th Int'l Colloquium on Automata, Languages and Programming (ICALP '03)*, LNCS 2719. Springer, 1068–1079.
- Mühlenbein, H. [1992]. How genetic algorithms really work: Mutation and hill-climbing. R. Männer and R. Manderick, editors, *Parallel Problem Solving from Nature 2 (PPSN)*. North-Holland, Amsterdam, 15–25.
- Nemirovsky, A. S. and Yudin, D. B. [1983]. *Problem Complexity and Method Efficiency in Optimization*. Wiley, New York.
- Neumann, F. and Wegener, I. [2004]. Randomized local search, evolutionary algorithms, and the minimum spanning tree problem. *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO)*, LNCS 3102. Springer, 713–724.
- Rechenberg, I. [1973]. *Evolutionsstrategie*. Frommann-Holzboog, Stuttgart, Germany.
- Rudolph, G. [1997]. *Convergence Properties of Evolutionary Algorithms*. Verlag Dr. Kovač, Hamburg.
- Scharnow, J., Tinnefeld, K., and Wegener, I. [2002]. Fitness landscapes based on sorting and shortest paths problems. *Parallel Problem Solving from Nature 7 (PPSN)*, LNCS 2439. Springer, 54–63.
- Schwefel, H.-P. [1995]. *Evolution and Optimum Seeking*. Wiley, New York.
- Wegener, I. [2001]. Theoretical aspects of evolutionary algorithms. *Proc. of the 28th Int'l Colloquium on Automata, Languages and Programming (ICALP)*, LNCS 2076. Springer, 64–78.