



# Quality gain analysis of the weighted recombination evolution strategy on general convex quadratic functions <sup>☆</sup>

Youhei Akimoto <sup>a,\*</sup>, Anne Auger <sup>b</sup>, Nikolaus Hansen <sup>b</sup>

<sup>a</sup> Faculty of Engineering, Information and Systems, University of Tsukuba, Japan

<sup>b</sup> Inria, RandOpt Team, CMAP, Ecole Polytechnique, France

## ARTICLE INFO

### Article history:

Received 15 November 2017

Received in revised form 9 April 2018

Accepted 10 May 2018

Available online 26 June 2018

### Keywords:

Evolution strategy

Weighted recombination

Quality gain analysis

Optimal step-size

General convex quadratic function

## ABSTRACT

Quality gain is the expected relative improvement of the function value in a single step of a search algorithm. Quality gain analysis reveals the dependencies of the quality gain on the parameters of a search algorithm, based on which one can derive the optimal values for the parameters. In this paper, we investigate evolution strategies with weighted recombination on general convex quadratic functions. We derive a bound for the quality gain and two limit expressions of the quality gain. From the limit expressions, we derive the optimal recombination weights and the optimal step-size, and find that the optimal recombination weights are independent of the Hessian of the objective function. Moreover, the dependencies of the optimal parameters on the dimension and the population size are revealed. Differently from previous works where the population size is implicitly assumed to be smaller than the dimension, our results cover the population size proportional to or greater than the dimension. Numerical simulation shows that the asymptotically optimal step-size well approximates the empirically optimal step-size for a finite dimensional convex quadratic function.

© 2018 Elsevier B.V. All rights reserved.

## 1. Introduction

**Background** Evolution Strategies (ES) are randomized search algorithms to minimize a black-box function  $f : \mathbb{R}^N \rightarrow \mathbb{R}$  in continuous domain, where neither the gradient nor the Hessian matrix of the objective function is available. The most advanced and commonly used category of evolution strategies is covariance matrix adaptation evolution strategy (CMA-ES) [2,3], which is recognized as the state-of-the-art black box continuous optimizer. It generates multiple candidate solutions from a multivariate normal distribution. They are evaluated on the objective function. The distribution parameters such as the mean vector and the covariance matrix are updated by using the candidate solutions and their ranking information, where the objective function values are not directly used. Due to its population-based and comparison-based nature, the algorithm is invariant to any strictly increasing transformation of the objective function in addition to the invariance to scaling, translation and rotation of the search space [4]. These invariance properties guarantee that the algorithm shows exactly the same behavior on a function  $f$  and on its transformation  $g \circ f \circ T$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is a strictly increasing function and  $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is a combination of scaling, translation and rotation defined as  $T : \mathbf{x} \mapsto a \cdot \mathbf{U}(\mathbf{x} - \mathbf{b})$  with a positive

<sup>☆</sup> This is the extension of our extended abstract presented at FOGA'2017 [1].

\* Corresponding author.

E-mail addresses: akimoto@cs.tsukuba.ac.jp (Y. Akimoto), anne.auger@inria.fr (A. Auger), nikolaus.hansen@inria.fr (N. Hansen).

real  $a > 0$ , an  $N$  dimensional vector  $\mathbf{b} \in \mathbb{R}^N$ , and an  $N$  dimensional orthogonal matrix  $\mathbf{U}$ . These invariance properties are the essence of the success of CMA-ES.

The performance evaluation of evolutionary algorithms is often based on empirical studies such as benchmarking on a test function suite [5,6] and well-considered performance assessment [7,8]. It is easier to check the performance of an algorithm on a specific problem in simulation than to analyze it mathematically. The invariance properties of an algorithm then generalize the empirical result to a class of infinitely many functions defined by the invariance relation. On the other hand, theoretical studies often require simplification of algorithms and assumptions on the objective function, because of the difficulty of the analysis of advanced algorithms due to their comparison-based and population-based nature and the complex adaptation mechanisms. Nevertheless, theoretical studies lead us to a better understanding of algorithms and reveals the dependency of the performance on the interval parameter settings. For example, the recombination weights in CMA-ES are selected based on the mathematical analysis of an evolution strategy [9].<sup>1</sup> The theoretical result of the optimal step-size on the sphere function is used to design a box constraint handling technique [10] and to design a termination criterion for a restart strategy [11]. A recent variant of CMA-ES [12] exploits the theoretical result of the optimal rate of convergence of the step-size to estimate the condition number of the product of the covariance matrix and the Hessian matrix of the objective function.

*Quality gain analysis* Quality gain and progress rate analysis [13–15] measure the expected progress of the mean vector in one step. On one side, differently from convergence analysis (e.g., [16]), analyses based on these quantities do not guarantee the convergence and often take a limit to derive an explicit formula. Moreover, the step-size adaptation and the covariance matrix adaptation are not taken into account. On the other side, one can derive quantitative explicit estimates of these quantities, which is not the case in convergence analysis. The quantitative explicit formulas are particularly useful to know the dependency of the expected progress on the parameters of the algorithm such as the population size, number of parents, and recombination weights, which we may not recognize from empirical studies of algorithms. The above mentioned recombination weights in CMA-ES are derived from the quality gain analysis of evolution strategies [9].

Although quality gain analysis is not meant to guarantee the convergence of the algorithm since it analyzes only a single step expected improvement, the progress rate is linked to the convergence rate of algorithms. It is directly related to the convergence rate of an “artificial” algorithm where the step-size is proportional to the distance to the optimum on the sphere function (see e.g., [17]). Moreover, the convergence rate of this artificial algorithm gives a bound on the convergence rate of algorithms that implement a proper step-size adaptation. For  $(1 + \lambda)$  or  $(1, \lambda)$  ESs the bound holds on any function with a unique global optimum; that is, any step-size adaptive  $(1 + \lambda)$ -ES optimizing any function  $f$  with a unique global optimum can not achieve a convergence rate faster than the convergence rate of the artificial algorithm on the sphere function where the step-size is the distance to the optimum times the optimal constant [18–20].<sup>2</sup> For algorithms implementing recombination, this bound still holds on spherical functions [19,20].

*Related work* In this paper, we investigate ESs with weighted recombination on a general convex quadratic function. ESs with weighted recombination samples multiple candidate solutions at one time and compute the weighted average of the candidate solutions to update the distribution mean vector. Weighted recombination ESs are among the most important categories of ESs since the standard CMA-ES and most of the recent variants of CMA-ES [21–23] employ weighted recombination.

The first analysis of weighted recombination ESs were done in [9], where the quality gain has been derived on the infinite dimensional sphere function  $f : x \mapsto \|x\|^2$ . The optimal step-size and the optimal recombination weights are derived. Reference [24] studied a variant of weighted recombination ESs called  $(\mu/\mu_1, \lambda)$ -ES, where  $(\mu/\mu_1, \lambda)$  stands for intermediate recombination, where the recombination weights are equal for the best  $\mu$  candidate solutions and zero for the other  $\lambda - \mu$  candidate solutions. The analysis has been performed on the quadratic functions with the Hessian  $\mathbf{A} = \frac{1}{2} \text{diag}(\alpha, \dots, \alpha, 1, \dots, 1)$ , where the number  $\lfloor N\theta \rfloor$  of diagonal elements that are  $\alpha > 1$  is controlled by the ratio  $\theta$  of short axes. Reference [25] studied the  $(1 + 1)$ -ES with the one-fifth success rule on the same function and showed the convergence rate of  $\Theta(1/(\alpha N))$ . Reference [26] studied ES with weighted recombination on the same function. Their results, progress rate and quality gain, depend on the so-called localization parameter, the steady-state value of which is then analyzed to obtain the steady-state quality gain. References [27,28] studied the progress rate and the quality gain of  $(\mu/\mu_1, \lambda)$ -ES on the general convex quadratic model.

The quality gain analysis and the progress rate analysis in the above listed references rely on a geometric intuition of the algorithm in the infinite dimensional search space and on various approximations. On the other hand, the rigorous derivation of the progress rate (or convergence rate of the algorithm with step-size proportional to the distance to the optimum) on the sphere function provided for instance in [17,20,29,30] only holds on spherical functions and provides solely a limit without a bound between the finite dimensional convergence rate and its asymptotic limit. The result of this

<sup>1</sup> The weights of CMA-ES were set before the publication [9] because the theoretical result of optimal weights on the sphere was known before the publication.

<sup>2</sup> More precisely,  $(1 + \lambda)$ -ES optimizing any function  $f$  (that may have more than one global optimum) can not converge towards a given optimum  $x^*$  faster in the search space than the artificial algorithm with step-size proportional to the distance to  $x^*$ .

paper is different in that we consider the general weighted recombination on the general convex quadratic objective and cover finite dimensional cases as well as the limit  $N \rightarrow \infty$ .

**Contributions** We study the weighted recombination ES on a general convex quadratic function  $f(x) = \frac{1}{2}(x - x^*)^T \mathbf{A}(x - x^*)$  on the finite  $N$  dimensional search space. We investigate the quality gain  $\phi$ , that is, the expectation of the relative function value decrease. We decompose  $\phi$  as the product of two functions:  $g$ , a function that depends only on the mean vector of the sampling distribution and the Hessian  $\mathbf{A}$ , and  $\bar{\phi}$ , the so-called *normalized quality gain* that depends essentially on all the algorithm parameters such as the recombination weights and the step-size. We approximate  $\bar{\phi}$  by an analytically tractable function  $\varphi$ . We call  $\varphi$  the *asymptotic normalized quality gain*. The main contributions are summarized as follows.

First, we derive the error bound between  $\bar{\phi}$  and  $\varphi$  for finite dimension  $N$ . To the best of our knowledge, this is the first work that performs the quality gain analysis for finite  $N$  and provides an error bound. The asymptotic normalized quality gain and the bounds in this paper are improved over the previous work [1]. Thanks to the explicit error bound derived in the paper, we can treat the population size  $\lambda$  increasing with  $N$  and provide (for instance) a rigorous sufficient condition on the dependency between  $\lambda$  and  $N$  such that the per-iteration quality gain scales with  $O(\lambda/N)$  for algorithms with intermediate recombination [15].

Second, we show that the error bound between  $\bar{\phi}$  and  $\varphi$  converges to zero as the learning rate  $c_m$  for the mean vector update tends to infinity. We derive the optimal step-size and the optimal recombination weights for  $\varphi$ , revealing the dependencies of these optimal parameters on  $\lambda$  and  $N$ . In contrast, the previous works of quality gain analysis mentioned above take the limit  $N \rightarrow \infty$  while  $\lambda$  is fix, hence assuming  $\lambda \ll N$ . Therefore, they do not reveal the dependencies of  $\bar{\phi}$  and the optimal parameters on  $\lambda$  when  $\lambda \ll N$ . We validate in experiments that the optimal step-size derived for  $c_m \rightarrow \infty$  provides a reasonable estimate of the optimal step-size even for  $c_m = 1$ .

Third, we prove that  $\varphi$  converges toward  $\bar{\phi}_\infty$  as  $N \rightarrow \infty$  under the condition  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 = 0$ , where  $\bar{\phi}_\infty$  is the limit of  $\bar{\phi}$  on the sphere function for  $N \rightarrow \infty$  derived in [9]. The condition  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 = 0$  holds, for example, for positive definite  $\mathbf{A}$  with bounded eigenvalues. It also holds for some positive semi-definite  $\mathbf{A}$  and for some positive definite  $\mathbf{A}$  with unbounded eigenvalues, for example with eigenvalues in  $[1, \sqrt{N}]$ . The result implies that the optimal recombination weights are independent of  $\mathbf{A}$ , whereas the optimal step-size heavily depends on  $\mathbf{A}$  and the distribution mean. This part of the contribution is a generalization of the previous foundation in [27,28], but the proof methodology is rather different. Furthermore, the error bound between  $\bar{\phi}$  and  $\varphi$  derived in this paper allows us to further investigate how fast  $\varphi$  converges toward  $\bar{\phi}_\infty$  as  $N \rightarrow \infty$ , depending on the eigenvalue distribution of  $\mathbf{A}$ .

**Organization** This paper is organized as follows. In Section 2, we formally define the evolution strategy with weighted recombination. The quality gain analysis on the infinite dimensional sphere function is revisited. In Section 3, we derive the quality gain bound for a finite dimensional convex quadratic function. In Section 4, important consequences of the quality gain bound are discussed. In Section 5, we conclude our paper. Properties of the normal order statistics that are important to understand our results are summarized in Appendix A and the detailed proofs of lemmas are provided in Appendix B.

**Notation** We apply the following mathematical notations throughout the paper. For integers  $n, m \in \mathbb{N}$  such that  $n \leq m$ , we denote the set of integers between  $n$  and  $m$  (including  $n$  and  $m$ ) by  $\llbracket n, m \rrbracket$ . Binomial coefficients are denoted as  $\binom{m}{n} = \frac{m!}{(m-n)!n!}$ . For real numbers  $a, b \in \mathbb{R}$  such that  $a \leq b$ , the open and the closed intervals are denoted as  $(a, b)$  and  $[a, b]$ , respectively. For an  $N$ -dimensional real vector  $x \in \mathbb{R}^N$ , let  $[x]_i$  denote the  $i$ -th coordinate of  $x$ . A sequence of length  $n$  is denoted as  $(x_i)_{i=1}^n = (x_1, \dots, x_n)$ , or just as  $(x_i)$ , and an infinite sequence is denoted as  $(x_i)_{i=1}^\infty$ . For  $x \in \mathbb{R}$ , the absolute value of  $x$  is denoted by  $|x|$ . For  $x \in \mathbb{R}^N$ , the Euclidean norm is denoted by  $\|x\| = (\sum_{i=1}^N [x]_i^2)^{\frac{1}{2}}$ . Let  $\mathbb{1}_{\text{condition}}$  be the indicator function which is 1 if *condition* is true and 0 otherwise. Let  $\Phi$  be the cumulative density function (c.d.f.) deduced by the (one-dimensional) standard normal distribution  $\mathcal{N}(0, 1)$ . Let  $\mathcal{N}_{i:\lambda}$  be the  $i$ -th smallest random variable among  $\lambda$  independently and standard normally distributed random variables, i.e.,  $\mathcal{N}_{1:\lambda} \leq \dots \leq \mathcal{N}_{\lambda:\lambda}$ . The expectation of a random variable (or vector)  $X$  is denoted as  $\mathbb{E}[X]$ . The conditional expectation of  $X$  given  $Y$  is denoted as  $\mathbb{E}[X | Y]$ . For a function  $f$  of  $n$  random variables  $(X_i)_{i=1}^n$ , the conditional expectation of  $F = f(X_1, \dots, X_n)$  given  $X_k$  for some  $k \in \llbracket 1, n \rrbracket$  is denoted as  $\mathbb{E}_k[F] = \mathbb{E}[F | X_k]$ . Similarly, the conditional expectation of  $F$  given  $X_k$  and  $X_l$  for different  $k, l \in \llbracket 1, n \rrbracket$  is denoted as  $\mathbb{E}_{k,l}[F] = \mathbb{E}[F | X_k, X_l]$ .

## 2. Formulation

### 2.1. Evolution strategy with weighted recombination

We consider an evolution strategy with weighted recombination. At each iteration  $t \geq 0$ , it draws  $\lambda$  independent random vectors  $Z_1, \dots, Z_\lambda$  from the  $N$ -dimensional standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , where  $\mathbf{0} \in \mathbb{R}^N$  is the zero vector and  $\mathbf{I}$  is the identity matrix of dimension  $N$ . The candidate solutions  $X_1, \dots, X_\lambda \sim \mathcal{N}(\mathbf{m}^{(t)}, (\sigma^{(t)})^2 \mathbf{I})$  are computed as  $X_i = \mathbf{m}^{(t)} + \sigma^{(t)} Z_i$ , where  $\mathbf{m}^{(t)} \in \mathbb{R}^N$  is the mean vector and  $\sigma^{(t)} > 0$  is the standard deviation, also called the step-size or the mutation strength. The candidate solutions are evaluated on a given objective function  $f: \mathbb{R}^N \rightarrow \mathbb{R}$ . Without loss of generality (w.l.o.g.), we assume  $f$  to be minimized. Let  $i:\lambda$  be the index of the  $i$ -th best candidate solution among

**Algorithm 1** Single step of the weighted recombination ES solving  $f$ .

```

1: procedure ES( $\mathbf{m}, \sigma, (w_k)_{k=1}^\lambda, c_m$ )
2:   for  $i = 1, \dots, \lambda$  do                                     ▷ Generate and evaluate  $\lambda$  candidate solutions
3:      $Z_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:      $X_i = \mathbf{m} + \sigma Z_i$ 
5:     Evaluate  $f(X_i)$ 
6:   end for
7:    $W(i; (X_k)_{k=1}^\lambda) = \sum_{k=1+l_i}^{u_i} w_k / (u_i - l_i)$                                      ▷ Compute the weights with (2)
8:    $\mathbf{m} \leftarrow \mathbf{m} + c_m \sigma \sum_{i=1}^\lambda W(i; (X_k)_{k=1}^\lambda) Z_i$                                      ▷ Update the mean with (3)
9:   return  $\mathbf{m}$ 
10: end procedure

```

$X_1, \dots, X_\lambda$ , i.e.,  $f(X_{1:\lambda}) \leq \dots \leq f(X_{\lambda:\lambda})$ , and  $w_1 \geq \dots \geq w_\lambda$  be the real-valued recombination weights. W.l.o.g., we assume  $\sum_{i=1}^\lambda |w_i| = 1$ . Let  $\mu_w = 1 / \sum_{i=1}^\lambda w_i^2$  denote the so-called effective variance selection mass. The mean vector is updated according to

$$\mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + c_m \sum_{i=1}^\lambda w_i (X_{i:\lambda} - \mathbf{m}^{(t)}) , \tag{1}$$

where  $c_m > 0$  is the learning rate of the mean vector update.

In this paper we reformulate (1) to investigate the algorithm with mathematical rigor. Hereunder, we write the candidate solutions,  $X_1, \dots, X_\lambda$ , and the corresponding random vectors,  $Z_1, \dots, Z_\lambda$ , as sequences  $(X_i)_{i=1}^\lambda$  and  $(Z_i)_{i=1}^\lambda$  for short. First, we introduce the weight function

$$W(i; (X_k)_{k=1}^\lambda) := \sum_{k=1+l_i}^{u_i} \frac{w_k}{u_i - l_i} , \quad \text{where} \quad \begin{cases} l_i = \sum_{j=1}^\lambda \mathbb{1}_{f(X_j) < f(X_i)} \\ u_i = \sum_{j=1}^\lambda \mathbb{1}_{f(X_j) \leq f(X_i)} \end{cases} , \tag{2}$$

i.e.,  $l_i$  and  $u_i$  are the numbers of strictly and weakly better candidate solutions than  $X_i$ , respectively. The weight value for  $X_i$  is the arithmetic average of the weights  $w_k$  for the tie candidate solutions. In other words, all the tie candidate solutions have the same weight values. If there is no tie, the weight value for the  $i$ -th best candidate solution  $X_{i:\lambda}$  is simply  $w_i$ . In the following, we drop the subscripts and the superscripts for sequences unless they are unclear from the context and write simply as  $(X_k) = (X_k)_{k=1}^\lambda$ . With the weight function, we rewrite the mean vector update (1) as

$$\mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + c_m \sum_{i=1}^\lambda W(i; (X_k)) (X_i - \mathbf{m}^{(t)}) . \tag{3}$$

The above update (3) is equivalent with the original update (1) if there is no tie among  $\lambda$  candidate solutions. If the objective function is a convex quadratic function, there will be no tie with probability one. Therefore, they are equivalent with probability one. Algorithm 1 summarizes the single step of the algorithm, where we rewrite (3) by using  $X_i - \mathbf{m}^{(t)} = \sigma^{(t)} Z_i$ .

The above formulation is motivated twofold. One is to well define the update even when there is tie. In our formulation, tie candidate solutions receive the equal recombination weights. The other is a technical reason. In (1) the already sorted candidate solutions  $X_{i:\lambda}$  are all correlated and they are not anymore normally distributed. However, they are assumed to be normally distributed in the previous work [9,27,28]. To ensure that such an approximation leads to the asymptotically true quality gain limit, a mathematically involved analysis has to be done. See [17,29,30] for details. In (3), the weight function explicitly includes the ranking computation and  $X_i$  are still independent and normally distributed. This allows us to derive the quality gain on a convex quadratic function rigorously.

2.2. Quality gain analysis on the spherical function

The quality gain is defined as the expectation of the relative decrease of the function value. Formally, it is the conditional expectation of the relative decrease of the function value conditioned on the mean vector  $\mathbf{m}^{(t)} = \mathbf{m}$  and the step-size  $\sigma^{(t)} = \sigma$ , defined as follows.

**Definition 1.** The quality gain of Algorithm 1 given  $\mathbf{m}^{(t)} = \mathbf{m}$  and  $\sigma^{(t)} = \sigma$  is

$$\phi(\mathbf{m}, \sigma) = \frac{\mathbb{E}[f(\mathbf{m}) - f(\text{ES}(\mathbf{m}, \sigma, (w_k)_{k=1}^\lambda, c_m))]}{f(\mathbf{m}) - f(x^*)} , \tag{4}$$

where  $x^* \in \mathbb{R}^N$  is (one of) the global minimum point of  $f$ . Note that the quality gain depends also on  $(w_k)_{k=1}^\lambda, c_m$ , and the dimension  $N$ .

**Results** Algorithm 1 solving a spherical function  $f(x) = \|x\|^2$  is analyzed in [9]. For this purpose, the *normalized step-size* and the *normalized quality gain* are introduced as

$$\bar{\sigma} = \sigma \frac{c_m N}{\|\mathbf{m}\|} \quad \text{and} \quad \bar{\phi}(\mathbf{m}, \bar{\sigma}) = \frac{N}{2} \phi\left(\mathbf{m}, \sigma = \frac{\bar{\sigma} \|\mathbf{m}\|}{c_m N}\right), \tag{5}$$

respectively. This normalization of the step-size suggests that  $\sigma$  is proportional to  $\|\mathbf{m}\|$  and inverse proportional to  $c_m$  and to  $N$ . The normalized step-size  $\bar{\sigma}$  is proportional to the ratio between the actual step-size and the distance between the current mean and the optimal solution. This reflects the scale invariance of the algorithm on the sphere function, that is, the single step response is solely determined by the normalized step-size. The dimension  $N$  in the numerator implies that the step-size  $\sigma$  needs to be inversely proportional to  $N$ . The normalized quality gain  $\bar{\phi}$  is simply the quality gain given  $\bar{\sigma}$  scaled by  $N/2$ . The scaling by  $N$  reflects that the convergence speed can not exceed  $O(1/N)$  for any comparison based algorithm [31]. By taking  $N \rightarrow \infty$ , the normalized quality gain converges pointwise (w.r.t.  $\bar{\sigma}$ ) to

$$\begin{aligned} \bar{\phi}_\infty(\bar{\sigma}, (w_k)) &:= \lim_{N \rightarrow \infty} \bar{\phi}(\mathbf{m}, \bar{\sigma}) = -\bar{\sigma} \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] - \frac{\bar{\sigma}^2}{2\mu_w} \\ &= \frac{\mu_w}{2} \left( \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] \right)^2 \left( 1 - \left( \frac{\bar{\sigma}}{\bar{\sigma}^*((w_k))} - 1 \right)^2 \right), \end{aligned} \tag{6}$$

where  $\bar{\sigma}^*((w_k))$  denotes the normalized step-size  $\bar{\sigma}$  optimizing  $\bar{\phi}_\infty$  given  $(w_k)$  and is given by

$$\bar{\sigma}^*((w_k)) = -\mu_w \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}]. \tag{7}$$

A formal proof of this result is presented in Theorem 2 of [29] relying on the uniform integrability of some random variable proved in [30].

Consider the optimal recombination weights that maximize  $\bar{\phi}_\infty$  in (6). The optimal recombination weights are given independently of  $\bar{\sigma}$  by

$$w_k^* = -\frac{\mathbb{E}[\mathcal{N}_{k:\lambda}]}{\sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]|} \tag{8}$$

and  $\bar{\phi}_\infty$  is written as

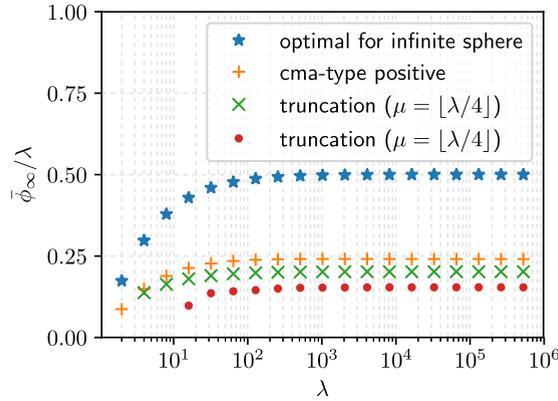
$$\bar{\phi}_\infty(\bar{\sigma}, (w_k^*)) = \frac{\sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2}{2} \left( 1 - \left( \frac{\bar{\sigma}}{\sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]|} - 1 \right)^2 \right). \tag{9}$$

Note that  $\bar{\sigma}^*((w_k^*)) = \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]|$ . Given  $\bar{\sigma}^*$  and  $(w_k^*)$ , we achieve the maximal value of  $\bar{\phi}_\infty$  that is  $\bar{\phi}_\infty(\bar{\sigma}^*((w_k^*)), (w_k^*)) = \sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2 / 2$ .

**Remarks** The optimal normalized step-size (7) and the normalized quality gain (6) given  $\bar{\sigma}^*$  depend on  $(w_k)$ . Particularly, they are proportional to  $\mu_w$ . For instance, under the optimal weights (8), we have  $-\sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] = \sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2 / \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]| \approx (\pi/2)^{1/2}$  for a sufficiently large  $\lambda$ .<sup>3</sup> Then, from (6) and (7) we know  $\bar{\sigma} \propto \mu_w$  and  $\bar{\phi}_\infty \propto \mu_w$ . Moreover, using the relation  $\mu_w = (\sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]|)^2 / \sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2 \approx (2/\pi)\lambda$ , we can reword it as that the optimal step-size and the normalized quality gain given  $\bar{\sigma}^*$  are proportional to  $\lambda$ . Fig. 1 shows how  $\bar{\phi}_\infty/\lambda$  scales with  $\lambda$  when the optimal step-size  $\sigma = \bar{\sigma}^*((w_k)) \|\mathbf{m}\| / (c_m N)$  is set. This shows that the normalized quality gain, and hence the optimal normalized step-size, are proportional to  $\lambda$  for standard weight schemes. When the optimal weights are used,  $\bar{\phi}_\infty/\lambda$  goes up to 0.5 as  $\lambda$  increases. On the other hand, nonnegative weights can not achieve the value of  $\bar{\phi}_\infty/\lambda$  above 0.25. The CMA type weights are designed to approximate the optimal nonnegative weights, where the first half are proportional to the optimal setting and the last half are zero. The truncation weights result in a smaller normalized quality gain. It is shown in [15] that the truncation weights achieve  $\bar{\phi}_\infty \in O(\mu \log(\lambda/\mu))$ .

The normalized quality gain limit  $\bar{\phi}_\infty$  depends only on the normalized step-size  $\bar{\sigma}$  and the weights  $(w_k)$ . Since the normalized step-size does not change if we multiply  $c_m$  by some factor and divide  $\sigma$  by the same factor,  $c_m$  does not have any impact on  $\bar{\phi}_\infty$ , hence on the quality gain  $\phi$ . This is unintuitive and is not true in a finite dimensional space. The step-size  $\sigma$  realizes the standard deviation of the sampling distribution and it has an impact on the ranking of the

<sup>3</sup> We used the facts  $\lim_{\lambda \rightarrow \infty} \sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2 / \lambda = 1$  and  $\lim_{\lambda \rightarrow \infty} \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]| / \lambda = (2/\pi)^{1/2}$ . See Appendix A for details.



**Fig. 1.** The normalized quality gain limit  $\bar{\phi}_\infty(\bar{\sigma}^*(w_k), (w_k))$  divided by  $\lambda$ . Four different weight schemes are employed: the optimal weights ( $w_k \propto -\mathbb{E}[\mathcal{N}_{k;\lambda}]$ ), the weights used in the CMA-ES ( $w_k \propto \max(\ln(\frac{\lambda+1}{2}) - \ln(k), 0)$ ), and the truncation weights ( $w_k = 1/\mu$  for  $k = 1, \dots, \mu$  and  $w_k = 0$  for  $k = \mu + 1, \dots, \lambda$ ) with  $\mu = \lfloor \lambda/4 \rfloor$  and  $\mu = \lfloor \lambda/10 \rfloor$ . All the weights are scaled so that  $\sum_{k=1}^{\lambda} |w_k| = 1$ . The value of  $\mathbb{E}[\mathcal{N}_{i;\lambda}]$  is approximated by the Blom’s formula (see Appendix A).

candidate solutions. On the other hand, the product  $\sigma c_m$  is the step-size of the  $\mathbf{m}$ -update that depends on the ranking of the candidate solutions. The normalized quality gain limit provided above tells us that the ranking of the candidate solutions is independent of  $\bar{\sigma}$  in the infinite dimensional space. We will discuss this further in Section 4.

The quality gain is to measure the improvement in one iteration. If we generate and evaluate  $\lambda$  candidate solutions every iteration, the quality gain per evaluation ( $f$ -call) is  $1/\lambda$  times smaller, i.e., the quality gain per evaluation is  $1/N$ , rather than  $\lambda/N$ . It implies that the number of iterations to achieve the same amount of the quality gain is inversely proportional to  $\lambda$ . This is the best we can hope for when the algorithm is implemented on a parallel computer. However, since the above result is obtained in the limit  $N \rightarrow \infty$  while  $\lambda$  is fixed, it is implicitly assumed that  $\lambda \ll N$ . The optimal down scaling of the number of iterations indeed only holds for  $\lambda \ll N$ . In practice, the quality gain per iteration tends to level out as  $\lambda$  increases. We will revisit this point in Section 4 and see how the optimal values for  $\bar{\sigma}$  and  $\bar{\phi}$  depend on  $N$  and  $\lambda$  when both are finite.

### 3. Quality gain analysis on general quadratic functions

In this section we investigate the normalized quality gain of Algorithm 1 minimizing a quadratic function with its Hessian  $\nabla \nabla f(x) = \mathbf{A}$  assumed to be nonnegative definite and symmetric, i.e.,

$$f(x) = \frac{1}{2}(x - x^*)^T \mathbf{A}(x - x^*) \quad (10)$$

where  $x^* \in \mathbb{R}^N$  is the global optimal solution.<sup>4</sup> W.l.o.g., we assume  $\text{Tr}(\mathbf{A}) = 1$ .<sup>5</sup> For the sake of notation simplicity we denote the directional vector of the gradient of  $f$  at  $\mathbf{m}$  by  $\mathbf{e} = \frac{\nabla f(\mathbf{m})}{\|\nabla f(\mathbf{m})\|} = \frac{\mathbf{A}(\mathbf{m} - x^*)}{\|\mathbf{A}(\mathbf{m} - x^*)\|}$ . To make the dependency of  $\mathbf{e}$  on  $\mathbf{m}$  clear, we sometimes write it as  $\mathbf{e}_m$ .

#### 3.1. Normalized quality gain and normalized step-size

We introduce the normalized step-size and the normalized quality gain. First of all, if the objective function is homogeneous around the optimal solution  $x^*$ , the optimal step-size must be a homogeneous function of degree 1 with respect to  $\mathbf{m} - x^*$ . This is formally stated in the following proposition. The proof is found in Appendix B.1.

**Proposition 2.** Let  $f : \mathbb{R}^N \rightarrow \mathbb{R}$  be a homogeneous function of degree  $n$ , i.e.,  $f(\alpha \cdot x) = \alpha^n f(x)$  for a fixed integer  $n > 0$  for any  $\alpha > 0$  and any  $x \in \mathbb{R}^N$ . Consider Algorithm 1 minimizing a function  $f^* : x \mapsto f(x - x^*)$ . Then, the quality gain is scale-invariant, i.e.,  $\phi(x^* + (\mathbf{m} - x^*), \sigma) = \phi(x^* + \alpha(\mathbf{m} - x^*), \alpha\sigma)$  for any  $\alpha > 0$ . Moreover, the optimal step-size  $\sigma^* = \text{argmax}_{\sigma \geq 0} \phi(\mathbf{m}, \sigma)$ , if it is well-defined, is a function of  $\mathbf{m} - x^*$ . For the sake of simplicity we write the optimal step-size as a map  $\sigma^* : \mathbf{m} - x^* \mapsto \sigma^*(\mathbf{m} - x^*)$ . It is a homogeneous function of degree 1, i.e.,  $\sigma^*(\alpha \cdot (\mathbf{m} - x^*)) = \alpha \sigma^*(\mathbf{m} - x^*)$  for any  $\alpha > 0$ .

<sup>4</sup> We use the following terminology in this paper. A nonnegative definite matrix  $\mathbf{A}$  is a matrix having only nonnegative eigenvalues, i.e.,  $x^T \mathbf{A} x \geq 0$  for all  $x \in \mathbb{R}^N$ . A nonnegative definite matrix  $\mathbf{A}$  is called positive definite if  $x^T \mathbf{A} x > 0$  for all  $x \in \mathbb{R}^N \setminus \{\mathbf{0}\}$ , otherwise it is called positive semi-definite. If  $\mathbf{A}$  is positive semi-definite, the optimum  $x^*$  is not unique.

<sup>5</sup> None of the algorithmic components and the quality measures used in the paper are affected by multiplying a positive constant to  $\mathbf{A}$ , or equivalently to  $f$ . To consider a general  $\mathbf{A}$ , simply replace  $\mathbf{A}$  with  $\mathbf{A}/\text{Tr}(\mathbf{A})$  in the following of the paper.

Note that the quadratic function is homogeneous of degree 2, and the function  $\mathbf{m} \mapsto \|\nabla f(\mathbf{m})\| = \|\mathbf{A}(\mathbf{m} - \mathbf{x}^*)\|$  is homogeneous of degree 1 around  $\mathbf{x}^*$ . The latter is our candidate for the optimal step-size. We define the normalized step-size, the scale-invariant step-size, and the normalized quality gain for a quadratic function as follows.

**Definition 3.** For a convex quadratic function (10), the *normalized step-size*  $\bar{\sigma}$  and the *scale-invariant step-size*  $\sigma$  given  $\bar{\sigma}$  are defined as  $\bar{\sigma} = (\sigma c_m) / \|\nabla f(\mathbf{m})\|$  and  $\sigma = (\bar{\sigma} / c_m) \|\nabla f(\mathbf{m})\|$ .

**Definition 4.** Let  $g : \mathbb{R}^N \rightarrow \mathbb{R}$  be the  $\mathbf{m}$ -dependent scaling factor of the normalized quality gain defined as  $g(\mathbf{m}) = \|\nabla f(\mathbf{m})\|^2 / f(\mathbf{m})$ . The *normalized quality gain* for a quadratic function is defined as  $\bar{\phi}(\mathbf{m}, \bar{\sigma}) = \phi(\mathbf{m}, \sigma = \bar{\sigma} \|\nabla f(\mathbf{m})\| / c_m) / g(\mathbf{m})$ .

Note that the normalized step-size and the normalized quality gain defined above coincide with (5) if  $f(x) = \|x\|^2 / (2N)$ , where  $\mathbf{A} = \mathbf{I}/N$ ,  $\nabla f(\mathbf{m}) = \mathbf{m}/N$  and  $g(\mathbf{m}) = 2/N$ . Moreover, they are equivalent to Eq. (4.104) in [15] introduced to analyze the  $(1 + \lambda)$ -ES and the  $(1, \lambda)$ -ES. The same normalized step-size has been used for  $(\mu/\mu_1, \lambda)$ -ES [27,28]. See Section 4.3.1 of [15] for the motivation of this normalization.

*Non-isotropic Gaussian sampling* Throughout the paper, we assume that the multivariate normal sampling distributions have an isotropic covariance matrix. We can generalize all the following results to an arbitrary positive definite symmetric covariance matrix  $\mathbf{C}$  by considering a linear transformation of the search space. Indeed, let  $f : x \mapsto \frac{1}{2}(x - \mathbf{x}^*)^\top \mathbf{A}(x - \mathbf{x}^*)$ , and consider the coordinate transformation  $x \mapsto y = \mathbf{C}^{-\frac{1}{2}}x$ . In the latter coordinate system the function  $f$  can be written as  $f(x) = \bar{f}(y) = \frac{1}{2}(y - \mathbf{C}^{-\frac{1}{2}}\mathbf{x}^*)^\top (\mathbf{C}^{\frac{1}{2}}\mathbf{A}\mathbf{C}^{\frac{1}{2}})(y - \mathbf{C}^{-\frac{1}{2}}\mathbf{x}^*)$ . The multivariate normal distribution  $\mathcal{N}(\mathbf{m}, \sigma^2\mathbf{C})$  is transformed into  $\mathcal{N}(\mathbf{C}^{-\frac{1}{2}}\mathbf{m}, \sigma^2\mathbf{I})$  by the same transformation. Then, it is easy to prove that the quality gain on the function  $f$  given the parameter  $(\mathbf{m}, \sigma, \mathbf{C})$  is equivalent to the quality gain on the function  $\bar{f}$  given  $(\mathbf{C}^{-\frac{1}{2}}\mathbf{m}, \sigma, \mathbf{I})$ . The normalization factor  $g(\mathbf{m})$  of the quality gain and the normalized step-size are then rewritten as

$$g(\mathbf{m}) = \frac{\|\mathbf{C}^{\frac{1}{2}}\mathbf{A}(\mathbf{m} - \mathbf{x}^*)\|^2}{f(\mathbf{m}) \text{Tr}(\mathbf{C}^{\frac{1}{2}}\mathbf{A}\mathbf{C}^{\frac{1}{2}})}, \quad \bar{\sigma} = \frac{\sigma c_m \text{Tr}(\mathbf{C}^{\frac{1}{2}}\mathbf{A}\mathbf{C}^{\frac{1}{2}})}{\|\mathbf{C}^{\frac{1}{2}}\mathbf{A}(\mathbf{m} - \mathbf{x}^*)\|}.$$

### 3.2. Conditional expectation of the weight function

The quadratic objective (10) can be written as

$$f(\mathbf{m} + \Delta) = f(\mathbf{m}) + \nabla f(\mathbf{m})^\top \Delta + \frac{1}{2} \Delta^\top \mathbf{A} \Delta. \tag{11}$$

The normalized quality gain on a convex quadratic function can be written as (using (11) with  $\Delta = \mathbf{m}^{(t+1)} - \mathbf{m}^{(t)}$  and substituting (3))

$$\begin{aligned} \bar{\phi}(\mathbf{m}, \bar{\sigma}) &= -\bar{\sigma} \sum_{i=1}^{\lambda} \mathbb{E}[\mathbb{E}_i[W(i; (X_k))] \mathbf{e}^\top Z_i] \\ &\quad - \frac{\bar{\sigma}^2}{2} \sum_{i=1}^{\lambda} \sum_{j=1}^{\lambda} \mathbb{E}[\mathbb{E}_{i,j}[W(i; (X_k))W(j; (X_k))] Z_i^\top \mathbf{A} Z_j], \end{aligned}$$

where  $X_k = \mathbf{m}^{(t)} + \sigma^{(t)}Z_k$ , and  $\mathbb{E}_i$  and  $\mathbb{E}_{i,j}$  are the conditional expectations given  $X_i$  and  $(X_i, X_j)$ , respectively.

The following lemma provides the expression of the conditional expectation of the weight function, which allows us to derive the bound for the difference between  $\bar{\phi}$  and  $\varphi$ . In the following, let

$$\begin{aligned} P_b(k; n, p) &= \binom{n}{k} p^k (1-p)^{n-k} \\ P_t(k, l; n, p, q) &= \binom{n}{l+k} \binom{l+k}{k} p^k q^l (1-(p+q))^{n-(k+l)} \end{aligned}$$

denote the probability mass functions of the binomial and trinomial distributions, respectively, where  $0 \leq k \leq n$ ,  $0 \leq l \leq n$ ,  $k+l \leq n$ ,  $0 \leq p \leq 1$ ,  $0 \leq q \leq 1$  and  $p+q \leq 1$ . The proof of the lemma is provided in Appendix B.2.

**Lemma 5.** Let  $X \sim \mathcal{N}(\mathbf{m}, \sigma^2\mathbf{I})$  and  $(X_i)_{i=1}^{\lambda}$  be  $\lambda$  i.i.d. copies of  $X$ . Let  $F_f(t) = \Pr[f(X) < t]$  be the c.d.f. of the function value  $f(X)$ . Then, we have for any  $i, j \in \llbracket 1, \lambda \rrbracket$ ,  $i \neq j$ ,

$$\begin{aligned} \mathbb{E}_i[W(i; (X_k))] &= u_1(F_f(f(X_i))) , \\ \mathbb{E}_i[W(i; (X_k))^2] &= u_2(F_f(f(X_i))) , \\ \mathbb{E}_{i,j}[W(i; (X_k))W(j; (X_k))] &= u_3(F_f(f(X_i)), F_f(f(X_j))) , \end{aligned}$$

where

$$u_1(p) = \sum_{k=1}^{\lambda} w_k P_b(k-1; \lambda-1, p) , \tag{12}$$

$$u_2(p) = \sum_{k=1}^{\lambda} w_k^2 P_b(k-1; \lambda-1, p) , \tag{13}$$

$$u_3(p, q) = \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l P_t(k-1, l-k-1; \lambda-2, \min(p, q), |q-p|) . \tag{14}$$

Thanks to Lemma 5 and the fact that  $(X_k)_{k=1}^{\lambda}$  are i.i.d., we can further rewrite the normalized quality gain as

$$\begin{aligned} \bar{\phi}(\mathbf{m}, \bar{\sigma}) &= -\bar{\sigma} \lambda \mathbb{E}[u_1(F_f(f(X))) \mathbf{e}^T Z] - \frac{\bar{\sigma}^2 \lambda}{2} \mathbb{E}[u_2(F_f(f(X)))(Z^T \mathbf{A} Z - 1)] \\ &\quad - \frac{\bar{\sigma}^2 \lambda}{2} \mathbb{E}[u_2(F_f(f(X)))] - \frac{\bar{\sigma}^2 (\lambda-1) \lambda}{2} \mathbb{E}[u_3(F_f(f(X)), F_f(f(\tilde{X}))) Z^T \mathbf{A} \tilde{Z}] . \end{aligned} \tag{15}$$

Here  $Z$  and  $\tilde{Z}$  are independent and  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ -distributed, and  $X = \mathbf{m} + \sigma Z$  and  $\tilde{X} = \mathbf{m} + \sigma \tilde{Z}$ , where  $\sigma = \bar{\sigma} \|\nabla f(\mathbf{m})\| / c_m$  is the scale-invariant step-size. Note that  $X$  and  $\tilde{X}$  are independent and  $\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{I})$ -distributed.

The following Lemma shows the Lipschitz continuity of  $u_1$ ,  $u_2$ , and  $u_3$ . The proof is provided in Appendix B.3.

**Lemma 6.** *The functions  $u_1$ ,  $u_2$ , and  $u_3$  are  $\ell_1$ -Lipschitz continuous, i.e.,  $|u_1(p_1) - u_1(p_2)| \leq L_1 |p_1 - p_2|$ ,  $|u_2(p_1) - u_2(p_2)| \leq L_2 |p_1 - p_2|$ , and  $|u_3(p_1, q_1) - u_3(p_2, q_2)| \leq L_3 (|p_1 - p_2| + |q_1 - q_2|)$ , with the Lipschitz constants*

$$\begin{aligned} L_1 &= \sup_{0 < p < 1} \left| (\lambda-1) \sum_{k=1}^{\lambda-1} (w_{k+1} - w_k) P_b(k-1; \lambda-2, p) \right| , \\ L_2 &= \sup_{0 < p < 1} \left| (\lambda-1) \sum_{k=1}^{\lambda-1} (w_{k+1}^2 - w_k^2) P_b(k-1; \lambda-2, p) \right| , \\ L_3 &= \max \left[ \sup_{0 < p < q < 1} \left| \sum_{k=1}^{\lambda-2} \sum_{l=k+2}^{\lambda} w_l (w_{k+1} - w_k) P_t(k-1, l-k-2; \lambda-3, p, q-p) \right| , \right. \\ &\quad \left. \sup_{0 < p < q < 1} \left| \sum_{k=1}^{\lambda-2} \sum_{l=k+2}^{\lambda} w_k (w_l - w_{l-1}) P_t(k-1, l-k-2; \lambda-3, p, q-p) \right| \right] (\lambda-2) . \end{aligned}$$

Upper bounds for the above Lipschitz constants are discussed in Appendix B.4.

### 3.3. Theorem: normalized quality gain on convex quadratic functions

The following main theorem provides the error bound between  $\bar{\phi}$  and  $\phi$ .

**Theorem 7.** *Consider Algorithm 1 and let  $f$  be a convex quadratic objective function (10). Let the normalized step-size  $\bar{\sigma}$  and the normalized quality gain  $\bar{\phi}$  defined in Definition 3 and Definition 4, respectively. Let  $\mathbf{e}_m = \nabla f(\mathbf{m}) / \|\nabla f(\mathbf{m})\|$  and  $\alpha = \min(1, (\bar{\sigma} / c_m) \text{Tr}(\mathbf{A}^2)^{1/2})$ . Define*

$$G(\alpha) = \min \left[ 1, \alpha \left( 2 + \frac{2^{\frac{1}{2}} (\ln(1/\alpha))^{\frac{1}{2}}}{\pi^{\frac{1}{2}}} + \frac{d_1(\mathbf{A}) \ln(1/\alpha)}{(2\pi)^{\frac{1}{2}} \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}} \right) \right] \tag{16}$$

and

$$\begin{aligned} \varphi(\bar{\sigma}, (w_k), \mathbf{e}_m, \mathbf{A}) &= -\bar{\sigma} \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] - \frac{\bar{\sigma}^2}{2} \sum_{i=1}^{\lambda} w_i^2 (1 - \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m) \\ &\quad - \frac{\bar{\sigma}^2}{2} \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m \sum_{i=1}^{\lambda} \sum_{j=1}^{\lambda} w_i w_j \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] , \end{aligned} \quad (17)$$

and let  $L_1, L_2, L_3$  be the Lipschitz constants of  $u_1, u_2$  and  $u_3$  defined in Lemma 5, respectively. Then,

$$\begin{aligned} \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} |\bar{\phi}(\mathbf{m}, \bar{\sigma}) - \varphi(\bar{\sigma}, (w_k), \mathbf{e}_m, \mathbf{A})| &\leq \bar{\sigma} \lambda L_1 ((2/\pi)^{\frac{1}{2}} G(\alpha) + (4\pi)^{-\frac{1}{2}} \alpha) \\ &\quad + \bar{\sigma} c_m \lambda L_2 (2^{-\frac{1}{2}} G(\alpha) + (8\pi)^{-\frac{1}{2}} \alpha) \alpha + \bar{\sigma} c_m \lambda (\lambda - 1) L_3 ((2/\pi)^{\frac{1}{2}} G(\alpha) + (2\pi^2)^{-\frac{1}{2}} \alpha) \alpha . \end{aligned} \quad (18)$$

The above theorem claims that if the right-hand side (RHS) of (18) is sufficiently small, the normalized quality gain  $\bar{\phi}$  is approximated by the asymptotic normalized quality gain  $\varphi$  defined in (17). Compared to  $\bar{\phi}_\infty$  in (6) derived for the infinite dimensional sphere function,  $\varphi$  is different even when  $\mathbf{A} \propto \mathbf{I}$ . We investigate the properties of  $\varphi$  in Section 4.1. The situations where the RHS of (18) is sufficiently small are discussed in Section 4.2 and Section 4.3. We remark that Theorem 3.4 in [1] provides a bound for the difference between  $\bar{\phi}$  and  $\bar{\phi}_\infty$ , instead of the difference between  $\bar{\phi}$  and  $\varphi$ . Introducing  $\varphi$  allows us to consider a finite dimensional case and to derive a tighter bound.

### 3.4. Outline of the proof of the main theorem

In the following of the section and in Appendix B, let  $Z_e = \mathbf{e}^T Z$ ,  $Z_\perp = Z - Z_e \mathbf{e}$ , and  $X = \mathbf{m} + \sigma Z$  for  $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Then,  $Z_e \sim \mathcal{N}(0, 1)$  and  $Z_\perp \sim \mathcal{N}(\mathbf{0}, \mathbf{I} - \mathbf{e} \mathbf{e}^T)$  and they are independent. Define

$$H_N = \frac{f(\mathbf{m} + \sigma Z) - \mathbb{E}[f(\mathbf{m} + \sigma Z)]}{\sigma \|\nabla f(\mathbf{m})\|} \quad \text{and} \quad h(Z) = \frac{1}{2} \frac{\bar{\sigma}}{c_m} (Z^T \mathbf{A} Z - 1) ,$$

where  $\mathbb{E}[f(\mathbf{m} + \sigma Z)] = f(\mathbf{m}) + \sigma^2/2$ . It is easy to see that  $H_N = Z_e + h(Z)$ . Let  $F_f$  and  $F_N$  be the c.d.f. induced by  $f(X)$  and  $H_N$ , respectively. Then,  $F_f(f(X)) = F_N(H_N)$ . Let  $\tilde{Z}$  be the i.i.d. copy of  $Z$  and define  $\tilde{Z}_e, \tilde{Z}_\perp, \tilde{X}$ , and  $\tilde{H}_N$  analogously.

The first lemma allows us to approximate  $F_N$  (hence  $F_f$ ) with the c.d.f.  $\Phi$  of the standard normal distribution. The proof is based on the Lipschitz continuity of  $\Phi$  and the tail bound of  $h(Z)$  proved in Lemma 1 of [32]. The detail is provided in Appendix B.5.

**Lemma 8.** Let  $\alpha$  and  $G(\alpha)$  be defined in Theorem 7. Then,  $\sup_{t \in \mathbb{R}} |F_N(t) - \Phi(t)| \leq G(\alpha)$ .

The following three lemmas are used to bound each term on the RHS of (15). The proofs are straight-forward from the Lipschitz continuity of  $u_1, u_2$  and  $u_3$  and Lemma 8. The detailed proofs are found in Appendix B.6, Appendix B.7, and Appendix B.8, respectively.

**Lemma 9.** Let  $L_1, \alpha$ , and  $G(\alpha)$  be the quantities appeared in Lemma 6 and Theorem 7. Then,

$$|\mathbb{E}[u_1(F_f(f(X)))Z_e] - \mathbb{E}[u_1(\Phi(Z_e))Z_e]| \leq L_1 ((2/\pi)^{\frac{1}{2}} G(\alpha) + (4\pi)^{-\frac{1}{2}} \alpha) .$$

**Lemma 10.** Let  $L_2, \alpha$ , and  $G(\alpha)$  be the quantities appeared in Lemma 6 and Theorem 7. Then,

$$\begin{aligned} &|\mathbb{E}[u_2(F_f(f(X)))(Z^T \mathbf{A} Z - 1)] - \mathbb{E}[u_2(\Phi(Z_e))(Z^T \mathbf{A} Z - 1)]| \\ &\leq L_2 (2^{\frac{1}{2}} G(\alpha) + (2\pi)^{-\frac{1}{2}} \alpha) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} . \end{aligned}$$

**Lemma 11.** Let  $L_3, \alpha$ , and  $G(\alpha)$  be the quantities appeared in Lemma 6 and Theorem 7. Then,

$$\begin{aligned} &|\mathbb{E}[u_3(F_f(f(X)), F_f(f(\tilde{X})))Z^T \mathbf{A} \tilde{Z}] - \mathbb{E}[u_3(\Phi(Z_e), \Phi(\tilde{Z}_e))Z^T \mathbf{A} \tilde{Z}]| \\ &\leq L_3 ((8/\pi)^{\frac{1}{2}} G(\alpha) + (2^{\frac{1}{2}}/\pi) \alpha) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} . \end{aligned}$$

The asymptotic normalized quality gain  $\varphi$  in (17) is obtained by replacing the c.d.f.  $F_f$  of  $f(X)$  in (15) with the c.d.f.  $\Phi$  of the standard normal distribution. The above lemmas are used to bound the difference between  $\bar{\phi}$  and  $\varphi$ . The following lemma provides the explicit form of each term of (17). The proof is straight-forward in light of Lemma 5. The detail can be found in Appendix B.9.

**Lemma 12.** The functions  $u_1, u_2$  and  $u_3$  defined in Lemma 5 satisfy the following properties:

$$\lambda \mathbb{E}[u_1(\Phi(Z_e))Z_e] = \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i;\lambda}] , \tag{19}$$

$$\lambda \mathbb{E}[u_2(\Phi(Z_e))] = \sum_{i=1}^{\lambda} w_i^2 , \tag{20}$$

$$\lambda \mathbb{E}[u_2(\Phi(Z_e))(Z^T \mathbf{A} Z - 1)] = \mathbf{e}^T \mathbf{A} \mathbf{e} \sum_{i=1}^{\lambda} w_i^2 (\mathbb{E}[\mathcal{N}_{i;\lambda}^2] - 1) , \tag{21}$$

$$\lambda(\lambda - 1) \mathbb{E}[u_3(\Phi(Z_e), \Phi(\tilde{Z}_e))Z^T \mathbf{A} \tilde{Z}] = 2 \mathbf{e}^T \mathbf{A} \mathbf{e} \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \mathbb{E}[\mathcal{N}_{k;\lambda} \mathcal{N}_{l;\lambda}] . \tag{22}$$

Now we finalize the proof of the main theorem. Using Lemma 12, we can rewrite (17) as

$$\begin{aligned} \varphi(\bar{\sigma}, (w_k)_{k=1}^{\lambda}, \mathbf{e}_m, \mathbf{A}) &= -\bar{\sigma} \lambda \mathbb{E}[u_1(\Phi(Z_e))Z_e] - \frac{\bar{\sigma}^2}{2} \lambda \mathbb{E}[u_2(\Phi(Z_e))] \\ &\quad - \frac{\bar{\sigma}^2}{2} \lambda \mathbb{E}[u_2(\Phi(Z_e))(Z^T \mathbf{A} Z - 1)] - \frac{\bar{\sigma}^2}{2} \lambda(\lambda - 1) \mathbb{E}[u_3(\Phi(Z_e), \Phi(\tilde{Z}_e))Z^T \mathbf{A} \tilde{Z}] . \end{aligned}$$

From the equation (15) and the above expression of  $\varphi$ , we have

$$\begin{aligned} \bar{\phi}(\mathbf{m}, \bar{\sigma}) - \varphi(\bar{\sigma}, (w_k)_{k=1}^{\lambda}, \mathbf{e}_m, \mathbf{A}) &= -\bar{\sigma} \lambda \mathbb{E}[(u_1(F_f(f(X))) - u_1(\Phi(Z_e)))Z_e] \\ &\quad - \frac{\bar{\sigma}^2 \lambda}{2} \mathbb{E}[(u_2(F_f(f(X))) - u_2(\Phi(Z_e)))] \\ &\quad - \frac{\bar{\sigma}^2 \lambda}{2} \mathbb{E}[(u_2(F_f(f(X))) - u_2(\Phi(Z_e)))(Z^T \mathbf{A} Z - 1)] \\ &\quad - \frac{\bar{\sigma}^2 (\lambda - 1) \lambda}{2} \mathbb{E}[(u_3(F_f(f(X)), F_f(f(\tilde{X}))) - u_3(\Phi(Z_e), \Phi(\tilde{Z}_e)))Z^T \mathbf{A} \tilde{Z}] . \end{aligned}$$

From the well-known fact (e.g., Theorem 2.1 of [33]) that for a random variable  $X$  with a continuous c.d.f.  $F_X$  the random variable  $F_X(X)$  is uniformly distributed on  $[0, 1]$ , we can prove both  $F_f(f(X))$  and  $\Phi(Z_e)$  are uniformly distributed on  $[0, 1]$ . Therefore, we have  $\mathbb{E}[u_2(F_f(f(X)))] = \mathbb{E}[u_2(\Phi(Z_e))] = \mathbb{E}[u_2(\mathcal{U}[0, 1])]$ , and the second term on the RHS of the above equality is zero. Applying the triangular inequality and Lemma 9, Lemma 10, and Lemma 11, we obtain (18). It completes the proof of Theorem 7.

#### 4. Consequences

Theorem 7 tells that if the RHS of (18) is sufficiently small, the normalized quality gain  $\bar{\phi}(\mathbf{m}, \bar{\sigma})$  is well approximated by  $\varphi(\bar{\sigma}, (w_k)_{k=1}^{\lambda}, \mathbf{e}_m, \mathbf{A})$  defined in (17). First we investigate the parameter values that are optimal for  $\varphi$ . Then, we consider the situations when the RHS of (18) is sufficiently small.

Let  $\mathbf{n}_{(\lambda)}$  be the  $\lambda$  dimensional column vector whose  $i$ -th component is  $\mathbb{E}[\mathcal{N}_{i;\lambda}]$  and  $\mathbf{N}_{(\lambda)}$  be the  $\lambda$  dimensional symmetric matrix whose  $(i, j)$ -th elements are  $\mathbb{E}[\mathcal{N}_{i;\lambda} \mathcal{N}_{j;\lambda}]$ . Let  $\mathbf{w}$  and  $\bar{\mathbf{w}}$  be the  $\lambda$  dimensional column vector whose  $i$ -th element is  $w_i$  and  $\bar{\sigma} w_i$ , respectively. Now (17) can be written as

$$\varphi(\bar{\mathbf{w}}, \mathbf{e}, \mathbf{A}) = -\bar{\mathbf{w}}^T \mathbf{n}_{(\lambda)} - \frac{1}{2} (1 - \mathbf{e}^T \mathbf{A} \mathbf{e}) \bar{\mathbf{w}}^T \bar{\mathbf{w}} - \frac{1}{2} (\mathbf{e}^T \mathbf{A} \mathbf{e}) \bar{\mathbf{w}}^T \mathbf{N}_{(\lambda)} \bar{\mathbf{w}} . \tag{23}$$

In the following we use the following asymptotically true approximation for a sufficiently large  $\lambda$  (see (A.2) in Appendix A)

$$\frac{\bar{\mathbf{w}}^T \mathbf{N}_{(\lambda)} \bar{\mathbf{w}}}{\lambda \|\bar{\mathbf{w}}\|^2} \approx \frac{(\bar{\mathbf{w}}^T \mathbf{n}_{(\lambda)})^2}{\|\bar{\mathbf{w}}\|^2 \|\mathbf{n}_{(\lambda)}\|^2} \approx \frac{(\bar{\mathbf{w}}^T \mathbf{n}_{(\lambda)})^2}{\lambda \|\bar{\mathbf{w}}\|^2} . \tag{24}$$

By “for a sufficiently large  $\lambda$ ”, we mean for a  $\lambda$  large enough to approximate the left hand side (LHS) of (24) by the right-most side (RHS). For a sufficiently large  $\lambda$ , (23) is approximated by

$$\varphi(\bar{\mathbf{w}}, \mathbf{e}, \mathbf{A}) \approx -\bar{\mathbf{w}}^T \mathbf{n}_{(\lambda)} - \frac{1}{2} (1 - \mathbf{e}^T \mathbf{A} \mathbf{e}) \bar{\mathbf{w}}^T \bar{\mathbf{w}} - \frac{1}{2} (\mathbf{e}^T \mathbf{A} \mathbf{e}) (\bar{\mathbf{w}}^T \mathbf{n}_{(\lambda)})^2 . \tag{25}$$

4.1. Asymptotic normalized quality gain and optimal parameters

As we mentioned in the previous section,  $\varphi$  in (17) is different from the normalized quality gain limit  $\bar{\phi}_\infty$  in (6). Consider the sphere function  $\mathbf{A} = \mathbf{I}/N$ ; then, since  $\mathbf{e}^T \mathbf{A} \mathbf{e} = 1/N$  for any  $\mathbf{e}$  with  $\|\mathbf{e}\| = 1$ , we have

$$\varphi(\bar{\sigma}, (\mathbf{w}_k), \mathbf{e}, \mathbf{A}) = \bar{\phi}_\infty(\bar{\sigma}, (\mathbf{w}_k)_{k=1}^\lambda) + \frac{\bar{\sigma}^2}{2N} \sum_{i=1}^\lambda w_i^2 - \frac{\bar{\sigma}^2}{2N} \sum_{i=1}^\lambda \sum_{j=1}^\lambda w_i w_j \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] .$$

Note that the second and the third terms on the RHS are proportional to  $1/N$ . By taking the limit for  $N$  to infinity, we have  $\varphi = \bar{\phi}_\infty$ . Therefore, the second and third terms describe how the finite dimensional cases are different from the infinite dimensional case. In particular, the last term prevents the quality gain from scaling up proportionally to  $\lambda$  when  $\lambda \ll N$ .

*Optimal recombination weights* The recombination weights optimal for  $\varphi$  are provided in the following proposition.

**Proposition 13.** *The asymptotic normalized quality gain  $\varphi$  (17) is optimized when  $\bar{\mathbf{w}}$  is the solution to the following linear system of equations*

$$(\mathbf{I} + \mathbf{e}^T \mathbf{A} \mathbf{e} (\mathbf{N}_{(\lambda)} - \mathbf{I})) \bar{\mathbf{w}} = -\mathbf{n}_{(\lambda)} , \tag{26}$$

where  $\bar{\sigma}$  and  $w_i$  are uniquely determined using the condition  $\sum_{i=1}^\lambda |w_i| = 1$ . Then the optimal value of  $\varphi$  is  $-\frac{1}{2} \mathbf{n}_{(\lambda)}^T \bar{\mathbf{w}}^*$  where  $\bar{\mathbf{w}}^*$  is the solution to the linear system (26).

**Proof.** We obtain (26) by taking the derivative of (23)  $\bar{\mathbf{w}}$  with respect to  $\bar{\mathbf{w}}$ , and requiring  $\partial\varphi(\bar{\mathbf{w}}, \mathbf{e}, \mathbf{A})/\partial[\bar{\mathbf{w}}]_i = 0$ . This ends the proof.  $\square$

First, consider the limit for  $N \rightarrow \infty$  while  $\lambda$  is fixed. As long as the largest eigenvalue  $d_1(\mathbf{A})$  of  $\mathbf{A}$  converges to zero as  $N \rightarrow \infty$ , i.e.,  $\lim_{N \rightarrow \infty} d_1(\mathbf{A}) = 0$ , we have  $\mathbf{e}^T \mathbf{A} \mathbf{e} \rightarrow 0$  as  $N \rightarrow 0$ . Then, (26) reads  $\bar{\mathbf{w}} = -\mathbf{n}_{(\lambda)}$ . Therefore, we have the same optimal recombination weights as the ones derived for the infinite dimensional sphere function.

Next, consider a finite dimensional case. If  $\lambda$  is sufficiently large, the optimality condition (26) is approximated by

$$\left( \frac{1}{\lambda} \left( 1 - \frac{1}{\mathbf{e}^T \mathbf{A} \mathbf{e}} \right) \mathbf{I} + \frac{1}{\mathbf{e}^T \mathbf{A} \mathbf{e}} \frac{\mathbf{n}_{(\lambda)} \mathbf{n}_{(\lambda)}^T}{\|\mathbf{n}_{(\lambda)}\|^2} \right) \bar{\mathbf{w}} = -\frac{\mathbf{n}_{(\lambda)}}{\lambda} .$$

The solution to the above approximated condition is given by  $\bar{\mathbf{w}} \propto -\mathbf{n}_{(\lambda)}$  independently of  $\mathbf{A}$  and  $\mathbf{e}$ . It means, for a sufficiently large  $\lambda$ , the optimal recombination weights are approximated by the weights (8) optimal for the infinite dimensional sphere function.

*Optimal normalized step-size* The optimal  $\bar{\sigma}$  under a given  $(\mathbf{w}_k)_{k=1}^\lambda$  is provided in the following proposition.

**Proposition 14.** *Given  $\mathbf{w} = (w_1, \dots, w_\lambda)$ , the asymptotic normalized quality gain (17) is maximized when the normalized step-size  $\bar{\sigma}$  is*

$$\bar{\sigma}^* = \frac{-\sum_{i=1}^\lambda w_i \mathbb{E}[\mathcal{N}_{i:\lambda}]}{\sum_{i=1}^\lambda w_i^2 (1 - \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m) + \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m \sum_{i=1}^\lambda \sum_{j=1}^\lambda w_i w_j \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}]} , \tag{27}$$

then  $\varphi(\bar{\sigma}^*, (\mathbf{w}_k), \mathbf{e}, \mathbf{A}) = \frac{\bar{\sigma}^*}{2} (-\sum_{i=1}^\lambda w_i \mathbb{E}[\mathcal{N}_{i:\lambda}])$ .

**Proof.** It is a straight-forward consequence from differentiating (17) with respect to  $\bar{\sigma}$  and solving  $\partial\varphi/\partial\bar{\sigma} = 0$ .  $\square$

For a sufficiently large  $\lambda$  (see (24)), one can rewrite and approximate (27) as

$$\begin{aligned} \bar{\sigma}^* &= \frac{-\mathbf{w}^T \mathbf{n}_{(\lambda)}}{(1 - \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m) \|\mathbf{w}\|^2 + \mathbf{e}_m^T \mathbf{A} \mathbf{e}_m \mathbf{w}^T \mathbf{N}_{(\lambda)} \mathbf{w}} \\ &\approx \frac{(\mathbf{e}_m^T \mathbf{A} \mathbf{e}_m)^{-1} \mu_w (-\mathbf{w}^T \mathbf{n}_{(\lambda)})}{(\mathbf{e}_m^T \mathbf{A} \mathbf{e}_m)^{-1} - 1 + \mu_w (-\mathbf{w}^T \mathbf{n}_{(\lambda)})^2} . \end{aligned} \tag{28}$$

Note again that  $-\mathbf{w}^T \mathbf{n}_{(\lambda)} = -\sum_{i=1}^\lambda w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] \in O(1)$  for the optimal weights, CMA-type non-negative weights, and truncation weights with fixed truncation ratio. To provide a better insight, consider the case of the sphere function ( $\mathbf{A} = \mathbf{I}/N$ ). Then, the RMS of (28) reads

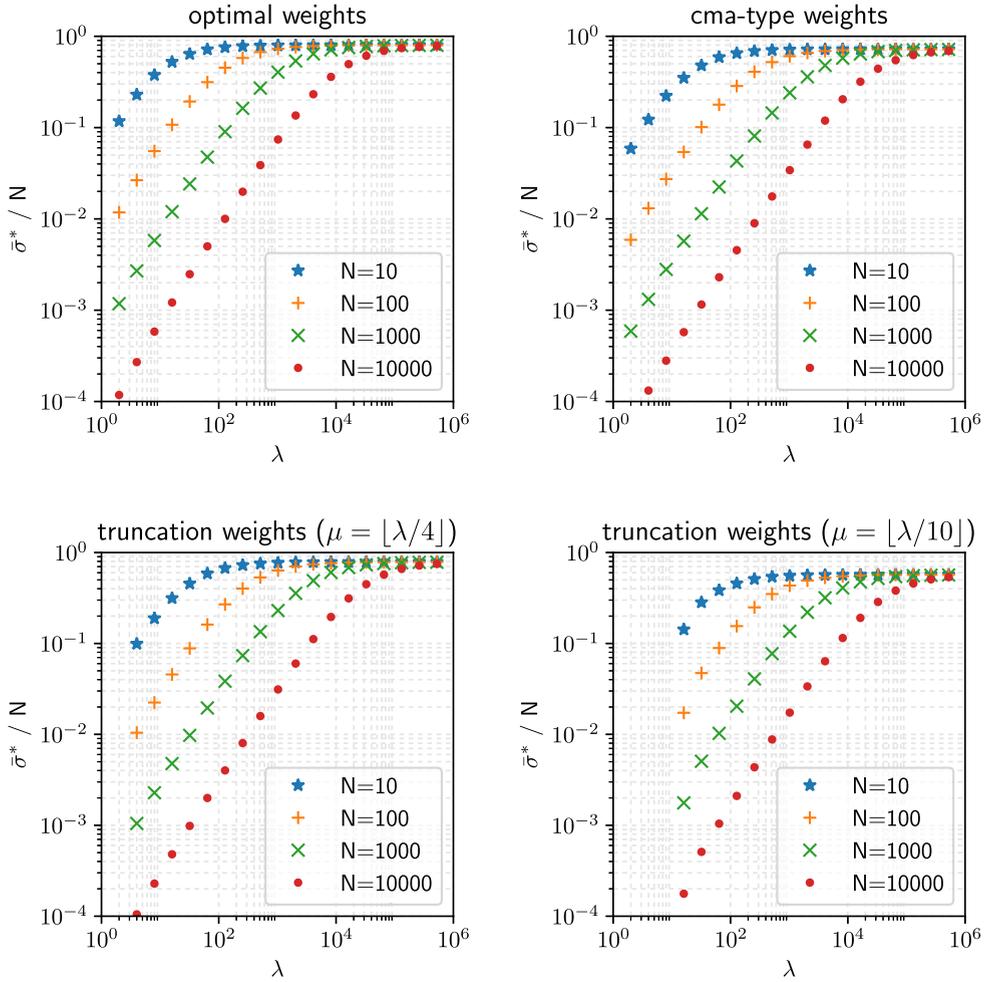


Fig. 2. Optimal normalized step-size on  $N = 10, 100, 1000, 10000$  dimensional sphere function for different weight schemes and different population size  $\lambda$ .

$$\bar{\sigma}^* \approx \frac{N\mu_w(-\mathbf{w}^T \mathbf{n}_{(\lambda)})}{N-1 + \mu_w(-\mathbf{w}^T \mathbf{n}_{(\lambda)})^2} \tag{29}$$

Then, we find the following: (i) if  $N \gg \mu_w$ , then  $\bar{\sigma}^* \approx \mu_w(-\mathbf{w}^T \mathbf{n}_{(\lambda)})$  and  $\varphi \approx \mu_w(-\mathbf{w}^T \mathbf{n}_{(\lambda)})^2/2$ ; (ii) if  $\mu_w \gg N$ , then  $\bar{\sigma}^* \approx N/(-\mathbf{w}^T \mathbf{n}_{(\lambda)})$  and  $\varphi \approx N/2$ . Fig. 2 visualizes the optimal normalized step-size (27) for various  $\mathbf{w}$  on the sphere function. The optimal normalized step-size (27) scales linearly for  $\lambda \leq N$  and it tends to level out for  $\lambda > N$ .

*Geometric interpretation of the optimal situation* On the infinite dimensional sphere function, we know that the optimal step-size puts the algorithm in the situation where  $f(\mathbf{m})$  improves twice as much by  $\mathbf{m}$  moving towards the optimum as it deteriorates by  $\mathbf{m}$  moving randomly in the subspace orthogonal to the gradient direction [13]. On a finite dimensional convex quadratic function, we find the analogous result. From (15) and lemmas in Section 3.4, the first term of the asymptotic normalized quality gain (17), i.e.  $-\bar{\sigma} \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i;\lambda}]$ , is due to the movement of  $\mathbf{m}$  in negative gradient direction, and the second and third terms are due to the random walk in the orthogonal subspaces.<sup>6</sup> The asymptotic normalized quality gain is maximized when the normalized step-size is set such that the first term is twice as large as the absolute value of the sum of the second and the third terms. That is, the amount of the decrease of  $f(\mathbf{m})$  by  $\mathbf{m}$  moving into the negative gradient direction is twice greater than the increase of  $f(\mathbf{m})$  by  $\mathbf{m}$  moving in its orthogonal subspace.

<sup>6</sup> More precisely, the second and the third terms come from the quadratic term in (11) that contains the information in the gradient direction as well. However, the above statement is true in the limit  $N \rightarrow \infty$  as long as  $\text{Tr}(\mathbf{A}^2) \rightarrow 0$ .

4.2. Infinitesimal step-size case

The RHS of (18), the error bound between  $\bar{\phi}$  and  $\varphi$ , converges to zero when  $\alpha \rightarrow 0$ . One such situation is the limit of  $\sigma / \|\mathbf{m}\| \rightarrow 0$  while  $c_m \sigma$  can remain positive, i.e., in the *mutate small, but inherit large* situation, which is formally stated in the next corollary.

**Corollary 15.** For any positive constant  $C > 0$ ,

$$\lim_{\sigma / \|\mathbf{m}\| \rightarrow 0} \sup_{\bar{\sigma} \in (0, C]} \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} \frac{1}{\bar{\sigma}} |\bar{\phi}(\mathbf{m}, \bar{\sigma}) - \varphi(\bar{\sigma}, (w_k), \mathbf{e}_m, \mathbf{A})| = 0 \tag{30}$$

**Proof.** Note that the function  $G(\alpha) \in \mathcal{O}(\alpha \ln(1/\alpha))$  as  $\alpha \rightarrow 0$ . Then, (18) reads

$$\begin{aligned} & \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} |\bar{\phi}(\mathbf{m}, \bar{\sigma}) - \varphi(\bar{\sigma}, (w_k), \mathbf{e}_m, \mathbf{A})| \\ & \in \bar{\sigma} \lambda \mathcal{O}(\alpha \ln(1/\alpha)) \left[ L_1 + c_m \alpha L_2 + c_m \alpha (\lambda - 1) L_3 \right]. \end{aligned} \tag{31}$$

Note also that

$$\alpha = \frac{\bar{\sigma}}{c_m} \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} = \frac{\sigma \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}}{\|\nabla f(\mathbf{m})\|} \leq \frac{\sigma}{\|\mathbf{m}\|} \frac{\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}}{d_N(\mathbf{A})}$$

and  $\alpha c_m = \bar{\sigma} \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} \leq C \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}$ . It implies that the RHS of (31) divided by  $\bar{\sigma}$  is in  $\mathcal{O}(\alpha \ln(1/\alpha)) \subseteq o(\alpha^{1-\epsilon})$  for any  $\epsilon > 0$  under the condition  $\bar{\sigma} \leq C$ . Since  $\alpha \rightarrow 0$  as  $\sigma / \|\mathbf{m}\| \rightarrow 0$ , (31) implies (30).  $\square$

If  $c_m$  is fixed, we have  $\bar{\sigma} \rightarrow 0$  as  $\sigma / \|\mathbf{m}\| \rightarrow 0$ . Then, the asymptotic normalized quality gain (17) converges towards zero as the bound on the RHS of (18) goes to zero. The above corollary tells that the bound converges faster than  $\bar{\sigma}$  does, while the asymptotic normalized quality gain decreases linearly in  $\bar{\sigma}$ . As a consequence, we find that the normalized quality gain approaches  $-\bar{\sigma} \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}]$  as  $\sigma / \|\mathbf{m}\| \rightarrow 0$ .

Consider the case that  $\bar{\sigma}$  is fixed, i.e.,  $c_m \sigma$  is fixed. Then, from the corollary we find that the normalized quality gain converges towards  $\varphi$  in (17) as  $\sigma / \|\mathbf{m}\| \rightarrow 0$ .

Taking  $c_m \rightarrow \infty$  we obtain  $\bar{\phi} \rightarrow \varphi$ . Though resembling a numerical gradient estimation when  $\sigma \rightarrow 0$ , it is not quite practical to take a large  $c_m$ . Indeed we usually rather do the opposite. For noisy optimization, the idea of *rescaled mutations* (corresponding to a large  $\sigma$  and a small  $c_m$ ) that was proposed by A. Ostermeier in 1993 (according to [13]) and analyzed in [34,35] is introduced to reduce the noise-to-signal ratio. If neither  $c_m \gg 1$  nor  $N \gg 1$ , the RHS of (18) will not be small enough to approximate the normalized quality gain by  $\varphi$  in (17). Then the normalized step-size defined in (27) is not guaranteed to provide an approximation of the optimal normalized step-size. However, in practice, we observe that the normalized step-size defined in (27) provides a reasonable approximation of the optimal normalized step-size for  $c_m \geq 1$ . We will see it in Fig. 4.

4.3. Infinite dimensional case

The other situation when  $\alpha \rightarrow 0$  occurs is the limit  $N \rightarrow \infty$  under the condition  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2) = 0$ . In this case, since  $\mathbf{e}^T \mathbf{A} \mathbf{e} \leq \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} \rightarrow 0$ , the limit expression  $\varphi$  converges to  $\bar{\phi}_{\infty}$ , i.e., the same limit on the sphere function. It is stated in the following corollary, which is a generalization of the result obtained in [9] from the sphere function to a general convex quadratic function.

**Corollary 16.** Let  $(\mathbf{A}_N)_{N=1}^{\infty}$  be the sequence of nonnegative definite matrices satisfying  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} = 0$ . Then,

$$\lim_{N \rightarrow \infty} \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} |\varphi(\bar{\sigma}, (w_k)_{k=1}^{\lambda}, \mathbf{e}, \mathbf{A}_N) - \bar{\phi}_{\infty}(\bar{\sigma}, (w_k)_{k=1}^{\lambda})| = 0, \tag{32}$$

where  $\bar{\phi}_{\infty}(\bar{\sigma}, (w_k)_{k=1}^{\lambda}) = -\bar{\sigma} \sum_{i=1}^{\lambda} w_i \mathbb{E}[\mathcal{N}_{i:\lambda}] - \bar{\sigma}^2 / (2\mu_w)$  as defined in (6). Moreover,

$$\lim_{N \rightarrow \infty} \sup_{\bar{\sigma} \in (0, C]} \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} |\bar{\phi}(\mathbf{m}, \bar{\sigma}) - \bar{\phi}_{\infty}(\bar{\sigma}, (w_k)_{k=1}^{\lambda})| = 0. \tag{33}$$

Corollary 16 shows that the normalized quality gain on a convex quadratic function converges towards the asymptotic normalized quality gain derived on the infinite dimensional sphere function as  $\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} \rightarrow 0$ . It implies that the optimal values of the recombination weights and the normalized step-size are independent of the Hessian of the objective function,

and given by (8). It is a nice feature since we do not need to tune the weight values depending on the function. Since any twice continuously differentiable function is locally approximated by a quadratic function, the optimal weights derived here are expected to be locally optimal on any twice continuously differentiable function.

In the above corollary, the population size  $\lambda$  is a constant over the dimension  $N$ . However, in the default setting of the CMA-ES, the population size is  $\lambda = 4 + \lfloor 3 \ln(N) \rfloor$ , meaning that the population size is unbounded. If the population size increases to infinity as  $N \rightarrow \infty$ , it is not guaranteed that the per-evaluation progress  $\bar{\phi}/\lambda$  converges to  $\bar{\phi}_\infty/\lambda$  as  $N \rightarrow \infty$ . The following proposition provides a sufficient condition on the recombination weights and the population size such that the per-evaluation progress  $\bar{\phi}/\lambda$  converges to  $\bar{\phi}_\infty/\lambda$  when  $\lambda$  increases as  $N$  increases.

**Proposition 17.** Let  $(\mathbf{A}_N)_{N=1}^\infty$  be the sequence of the Hessian matrix that satisfies  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}_N^2)^{\frac{1}{2}} = 0$ . Let  $(\lambda_N)_{N=1}^\infty$  be the sequence of the population size and  $(w_k^N)_{k=1}^{\lambda_N}$  be the sequence of the weights for the population size  $\lambda_N$ . Suppose for an arbitrarily small positive  $\epsilon$ ,

$$\lambda_N^2 \in o\left(\frac{1}{d_1(\mathbf{A}_N)}\right) \quad \text{and}$$

$$\max\left(\lambda_N, L_1^{\frac{1}{1-\epsilon}} \lambda_N^{\frac{2-\epsilon}{1-\epsilon}}, L_2^{\frac{1}{2-\epsilon}} \lambda_N^{\frac{3-\epsilon}{2-\epsilon}}, L_3^{\frac{1}{2-\epsilon}} \lambda_N^{\frac{4-\epsilon}{2-\epsilon}}\right) \in O\left(\frac{1}{\text{Tr}(\mathbf{A}_N^2)^{\frac{1}{2}}}\right).$$

Then,

$$\lim_{N \rightarrow \infty} \sup_{\bar{\sigma} \in (0, 2\bar{\sigma}^*)} \sup_{\mathbf{m} \in \mathbb{R}^N \setminus \{\mathbf{0}\}} \frac{1}{\lambda_N} \left| \bar{\phi}(\mathbf{m}, \bar{\sigma}) - \bar{\phi}_\infty(\bar{\sigma}, (w_k^N)_{k=1}^{\lambda_N}) \right| = 0, \tag{34}$$

where  $\bar{\sigma}^*$  is the normalized step-size optimal for  $\bar{\phi}_\infty(\bar{\sigma}, (w_k^N)_{k=1}^{\lambda_N})$ , which is formulated in (7).

**Proof.** A sufficient condition for  $\varphi$  to converge to  $\bar{\phi}_\infty$  for  $\bar{\sigma} \in (0, 2\bar{\sigma}^*)$  is that the third term on the RHS of (17) converges to zero as  $N \rightarrow \infty$ . As we know from Appendix A that  $\mathbf{w}^T \mathbf{N}_{(\lambda)} \mathbf{w} \leq d_1(\mathbf{N}_{(\lambda)}) \|\mathbf{w}\|^2 \leq \text{Tr}(\mathbf{N}_{(\lambda)}) \|\mathbf{w}\|^2 = \lambda/\mu_w$ . On the other hand,  $\mathbf{e}_m^T \mathbf{A} \mathbf{e}_m$  is no greater than the greatest eigen value  $d_1(\mathbf{A})$  of  $\mathbf{A}$ . The third term on the RHS of (17) is maximized when  $\bar{\sigma} = \bar{\sigma}^* = \mu_w(-\mathbf{w}^T \mathbf{n}_{(\lambda)})$ , where,  $\mu_w(\mathbf{w}^T \mathbf{n}_{(\lambda)}) \leq \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]|$  and  $\frac{1}{\lambda} \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]| \rightarrow (2/\pi)^{\frac{1}{2}}$ . From these arguments derives that the third term on the RHS of (17) converges to zero as  $N \rightarrow \infty$  provided that  $\lambda^2 d_1(\mathbf{A}) \rightarrow 0$  as  $N \rightarrow \infty$ .

Next we consider the convergence of the bound (RHS of (31)). Remember that  $\alpha = (\bar{\sigma}/c_m) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}$  and  $G(\alpha) \in O(\alpha \ln(1/\alpha))$ . For  $\bar{\sigma} \in (0, 2\bar{\sigma}^*)$ , we have  $(\bar{\sigma}/c_m) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} \leq 2(\bar{\sigma}^*/c_m) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}$ . Since  $\bar{\sigma}^* \in O(\lambda)$ , we have  $\bar{\sigma} \in O(\lambda)$  and  $\alpha \in O(\lambda \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}})$ . Then, the RHS of (31) divided by  $\lambda$ ,

$$\begin{aligned} & O(\bar{\sigma} \alpha \ln(1/\alpha) [L_1 + L_2 \alpha + L_3 \lambda \alpha]) \\ & \subseteq o(\bar{\sigma} \alpha^{1-\epsilon} [L_1 + L_2 \alpha + L_3 \lambda \alpha]) \\ & \subseteq o(\lambda^{2-\epsilon} \text{Tr}(\mathbf{A}^2)^{\frac{1-\epsilon}{2}} [L_1 + L_2 \lambda \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} + L_3 \lambda^2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}]), \end{aligned}$$

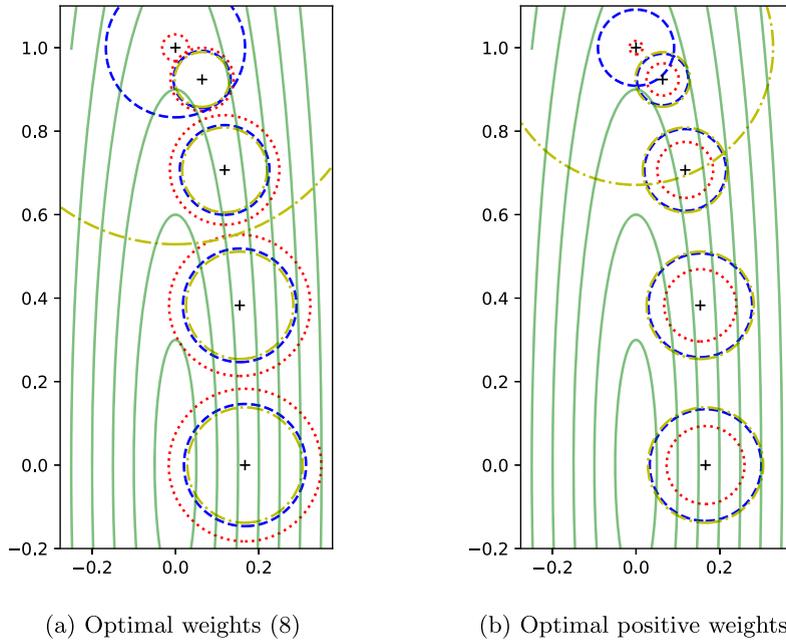
where the convergence of each term is supposed in the proposition.  $\square$

Consider the truncation weights with a fixed selection ratio  $\lambda = \rho\mu$  for some  $\rho > 1$  and the sequence of the Hessian matrices such that the condition number is bounded. From Appendix B.4, we have that  $L_1 \in O(\lambda^{-1/2})$ ,  $L_2 \in O(\lambda^{-3/2})$ , and  $L_3 \in O(\lambda^{-3/2})$ . Moreover, we have  $1/d_1(\mathbf{A}) \in O(N)$  and  $1/\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} \in O(N^{\frac{1}{2}})$ . Then, the condition of Proposition 17 reduces to  $\lambda \in o(N^{\frac{1}{3}})$ . This condition is a rigorous (but probably not tight) bound for the scaling of  $\lambda$  such that the per-iteration convergence rate of a  $(\mu/\mu, \lambda)$ -ES with a fixed  $\lambda/\mu$  on the sphere function scales like  $O(\lambda/N)$  [15, Equation 6.140]. One can also deduce the condition for the optimal weights and the CMA-type positive weights (positive half of the optimal weights).

#### 4.4. Effect of the eigenvalue distribution of the Hessian matrix

Corollary 16 implies that the optimal recombination weights are independent of the Hessian in the limit  $N \rightarrow \infty$  as long as  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2) = 0$ . Moreover, Proposition 13 and Corollary 15 together imply that the same values approximate the optimal recombination weights for a sufficiently large  $\lambda$  in the limit of  $\sigma/\|\mathbf{m}\| \rightarrow 0$ . On the other hand, the step-size and the progress rate depend on the Hessian. In the following we discuss the effect of the Hessian followed by a simulation. To make the discussion more intuitive, we remove the condition  $\text{Tr}(\mathbf{A}) = 1$  and consider an arbitrary non-negative definite symmetric  $\mathbf{A}$ . All the statements above still hold by replacing  $\mathbf{A}$  with  $\mathbf{A}/\text{Tr}(\mathbf{A})$ .

Given  $(w_k)_{k=1}^{\lambda}$ , the optimal normalized step-size  $\bar{\sigma}^*$  and the normalized quality gain  $\bar{\phi}_\infty(\bar{\sigma}, (w_k))$  are independent of the Hessian  $\mathbf{A}$  and the distribution mean  $\mathbf{m}$ . However, the step-size  $\sigma = (\bar{\sigma}/c_m) \|\nabla f(\mathbf{m})\|$  and the quality gain



**Fig. 3.** The asymptotically optimal step-size  $\sigma^* = \bar{\sigma}^* \|\nabla f(\mathbf{m})\|$  on  $f(x) = x^T \mathbf{A}x/2$  with  $\mathbf{A} = \text{diag}(1, 36)$ . The circles with radius  $\bar{\sigma}^* \|\nabla f(\mathbf{m})\|$  centered at  $\mathbf{m} = 2\mathbf{A}^{-\frac{1}{2}}(\cos(\theta), \sin(\theta))$  with  $\theta = \pi/2, 3\pi/8, \pi/4, \pi/8, 0$  are displayed, where the asymptotically optimal normalized step-size  $\bar{\sigma}^*$  is computed using (27) with the optimal weights (8) (left) and with the optimal positive weights (right). Red dotted:  $\lambda = 2$ , Blue dashed:  $\lambda = 10$ , Yellow dot-dashed:  $\lambda = 50$ . (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

$\phi(\mathbf{m}, \sigma) = g(\mathbf{m})\bar{\phi}_\infty(\bar{\sigma}, (w_k)_{k=1}^\lambda)$  depend on them through  $\|\nabla f(\mathbf{m})\|$  and  $g(\mathbf{m}) = \|\nabla f(\mathbf{m})\|^2 / f(\mathbf{m})$ . If  $\mathbf{m}$  is on a contour ellipsoid ( $f(\mathbf{m}) = 1$  for example),  $g(\mathbf{m})$  increases as  $\|\nabla f(\mathbf{m})\|$ . In other words, the greater the optimal step-size is, the greater quality gain we achieve. These quantities are bounded as

$$\frac{d_N(\mathbf{A})}{\text{Tr}(\mathbf{A})} \|\mathbf{m} - x^*\| \leq \|\nabla f(\mathbf{m})\| \leq \frac{d_1(\mathbf{A})}{\text{Tr}(\mathbf{A})} \|\mathbf{m} - x^*\| \text{ and}$$

$$\frac{d_N(\mathbf{A})}{\text{Tr}(\mathbf{A})} \leq \frac{g(\mathbf{m})}{2} \leq \frac{d_1(\mathbf{A})}{\text{Tr}(\mathbf{A})} .$$

The lower and upper equalities for both of the above inequalities hold if and only if  $\mathbf{m} - x^*$ , or equivalently  $\mathbf{e}_m$ , is parallel to the eigenspace corresponding to the smallest and largest eigenvalues of  $\mathbf{A}$ , respectively. Therefore, the optimal step-size and the quality gain can be different by the factor of at most  $\text{Cond}(\mathbf{A}) = d_1(\mathbf{A})/d_N(\mathbf{A})$ . Fig. 3 visualizes example cases. The asymptotic optimal step-size heavily depends on the location of  $\mathbf{m}$  if  $\mathbf{A}$  is ill-conditioned. If we focus on the area around each circle, the function landscape looks like a parabolic ridge function. Note that a relatively large step-size displayed at  $\mathbf{m} = (0, 1)$  for  $\lambda > 2$  is because  $\mathbf{e}_m^T \mathbf{A} \mathbf{e}_m \ll 1$  in (27), resulting in  $\bar{\sigma}^* \propto \mu_w$ . The asymptotic normalized quality gain is derived for the limit  $\sigma / \|\mathbf{m}\| \rightarrow 0$ , and the update of the mean vector results in an approximation of the negative gradient direction. If the mean vector is exactly on the longest axis of the hyper-ellipsoid, the gradient points to the optimal solution and a large normalized step-size is desired. However, this never happens in practice, since the mean vector will not be exactly in such a situation with probability one. We also remark that the asymptotically optimal normalized step-size (27) is not monotonically increasing w.r.t.  $\lambda$ . Indeed, we see in Fig. 3a smaller step-sizes for greater  $\lambda$  values, whereas they are monotonic in Fig. 3b. The main difference is that the step-sizes with the optimal weights for  $\lambda = 2$  can be greater than those with the optimal positive weights. Nonetheless, there is no guarantee that these figures reflect the *actually* optimal step-size precisely since displayed are the step-size optimal in the limit of  $c_m$  to infinity. Further investigation needs to be conducted.

Table 1 summarizes  $d_N(\mathbf{A})/\text{Tr}(\mathbf{A})$ ,  $d_1(\mathbf{A})/\text{Tr}(\mathbf{A})$  and  $\text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2$  for different types of  $\mathbf{A}$ . The greater the first two quantities are, the greater the optimal step-size and hence the quality gain are. The smaller the last quantity is, the more reliable it is to approximate  $\bar{\phi}$  with  $\bar{\phi}_\infty$ . If the condition number  $\alpha = \text{Cond}(\mathbf{A})$  is fixed, the worst case ( $d_N(\mathbf{A})/\text{Tr}(\mathbf{A})$ ) is maximized when the function has a discus type structure and is minimized when the function has a cigar type structure. The value of  $d_N(\mathbf{A})/\text{Tr}(\mathbf{A})$  will be close to  $1/N$  as  $N \rightarrow \infty$  for the discus type function, whereas it will be close to  $1/(N\alpha)$  for the cigar. Therefore, the discus type function is as easy to solve as the sphere function if  $N \gg \alpha$ , while the cigar type function takes roughly  $1/\alpha$  times more iterations to reach the same target function value. On the other hand, the inequality  $\text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 < 1/(N - 1)$  holds independently of  $\alpha$  on the cigar type function, while  $\text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2$  depends heavily on  $\alpha$

**Table 1**

Different types of the eigenvalue distributions of  $\mathbf{A}$ . The second to fourth types (discus:  $d_1(\mathbf{A}) = \alpha$  and  $d_2(\mathbf{A}) = \dots = d_N(\mathbf{A}) = 1$ , ellipsoid:  $d_i(\mathbf{A}) = \alpha^{\frac{i-1}{N-1}}$ , cigar:  $d_1(\mathbf{A}) = \dots = d_{N-1}(\mathbf{A}) = \alpha$  and  $d_N(\mathbf{A}) = 1$ ) have the condition number  $\text{Cond}(\mathbf{A}) = d_1(\mathbf{A})/d_N(\mathbf{A}) = \alpha$ , while the last type has the condition number  $N$ .

Type	$\frac{d_N(\mathbf{A})}{\text{Tr}(\mathbf{A})}$	$\frac{d_1(\mathbf{A})}{\text{Tr}(\mathbf{A})}$	$\frac{\text{Tr}(\mathbf{A}^2)}{\text{Tr}(\mathbf{A})^2}$
Sphere	$\frac{1}{N}$	$\frac{1}{N}$	$\frac{1}{N}$
Discus	$\frac{1}{(N-1)+\alpha}$	$\frac{\alpha}{(N-1)+\alpha}$	$\frac{(N-1)+\alpha^2}{((N-1)+\alpha)^2}$
Ellipsoid	$\frac{\alpha^{\frac{1}{N-1}}-1}{\alpha^{\frac{N}{N-1}}-1}$	$\frac{\alpha^{\frac{N}{N-1}}-\alpha}{\alpha^{\frac{N}{N-1}}-1}$	$\frac{(\alpha^{\frac{2N}{N-1}}-1)/(\alpha^{\frac{2}{N-1}}-1)}{(\alpha^{\frac{N}{N-1}}-1)^2/(\alpha^{\frac{1}{N-1}}-1)^2}$
Cigar	$\frac{1}{(N-1)\alpha+1}$	$\frac{\alpha}{(N-1)\alpha+1}$	$\frac{(N-1)\alpha^2+1}{((N-1)\alpha+1)^2}$
$d_i(\mathbf{A}) = i$	$\frac{1}{N(N+1)/2}$	$\frac{1}{(N+1)/2}$	$\frac{\frac{1}{6}N(N+1)(2N+1)}{(N(N+1)/2)^2}$

on the discus type function. The fraction will not be sufficiently small and we can not approximate the normalized quality gain by  $\bar{\phi}_\infty$  unless  $\alpha \ll N$  holds.<sup>7</sup>

The condition  $\lim_{N \rightarrow \infty} \text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 = 0$  also hold for some positive semi-definite  $\mathbf{A}$ , where only  $M < N$  eigenvalues of  $\mathbf{A}$  are positive and the others are zero. That is,  $d_1(\mathbf{A}) \geq \dots \geq d_M(\mathbf{A}) > 0$  and  $d_{M+1}(\mathbf{A}) = \dots = d_N(\mathbf{A})$ . In this case, the condition  $\text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 \rightarrow 0$  holds only if the dimension  $M$  of the effective search space tends to infinity as  $N \rightarrow \infty$ . The above inequalities are refined as follows. Let  $\mathbf{m}^+$  and  $\mathbf{m}^-$  be the decomposition of  $\mathbf{m}$  such that  $\mathbf{m}^-$  is the projection of  $\mathbf{m}$  onto the hyper-plane through  $x^*$  spanned by the eigenvectors of  $\mathbf{A}$  corresponding to the zero eigenvalue, and  $\mathbf{m}^+ = \mathbf{m} - \mathbf{m}^-$ . Then,

$$\frac{d_M(\mathbf{A})}{\text{Tr}(\mathbf{A})} \leq \frac{g(\mathbf{m})}{2} \leq \frac{d_1(\mathbf{A})}{\text{Tr}(\mathbf{A})}$$

$$\frac{d_M(\mathbf{A})}{\text{Tr}(\mathbf{A})} \|\mathbf{m}^+\| \leq \frac{\|\nabla f(\mathbf{m})\|}{\text{Tr}(\mathbf{A})} \leq \frac{d_1(\mathbf{A})}{\text{Tr}(\mathbf{A})} \|\mathbf{m}^+\| .$$

In this case,  $g(\mathbf{m})$  can be  $2/M$  if  $d_1(\mathbf{A}) = \dots = d_M(\mathbf{A}) > 0$  and  $d_i(\mathbf{A}) = 0$  for  $i \in \llbracket M + 1, N \rrbracket$ . The quality gain is then proportional to  $2/M$ , instead of  $2/N$ . That is, the evolution strategy with the optimal step-size solves the quadratic function with the effective rank  $M$  defined on the  $N$  dimensional search space as efficiently as it solves its projection onto the effective search space.

*Comment on the algorithm dynamics* The asymptotic quality gain depends on the distribution mean  $\mathbf{m}$  through  $g(\mathbf{m})$ . In practice, we observe near worst case performance with  $g(\mathbf{m}) \approx 2d_N(\mathbf{A})/\text{Tr}(\mathbf{A})$ , which implies that  $\mathbf{m} - x^*$  is almost parallel to the eigenspace corresponding to the smallest eigenvalue  $d_N(\mathbf{A})$  of the Hessian matrix. We provide an intuition to explain this behavior, which will be useful to understand the algorithm, even though the argument is not fully rigorous.

Consider Algorithm 1 with scale-invariant step-size (Definition 3). Lemma 9 implies that the order of the function values  $f(X_i)$  coincide with the order of  $[\mathcal{N}_i]_1 = \mathbf{e}^T(X_i - \mathbf{m}^{(t)})/\sigma^{(t)}$ , where  $\mathbf{e} = \nabla f(\mathbf{m}^{(t)})/\|\nabla f(\mathbf{m}^{(t)})\|$ . This is because if  $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , then  $Z^T \mathbf{A} Z/\text{Tr}(\mathbf{A})$  in (11) almost surely converges to one by the strong law of large numbers as  $N \rightarrow \infty$  under  $\text{Tr}(\mathbf{A}^2)/\text{Tr}(\mathbf{A})^2 \rightarrow 0$ . It means that the function value of a candidate solution is determined solely by the first component on the right-hand side of (11), that is,  $\mathbf{e}^T(X_i - \mathbf{m}^{(t)})/\sigma^{(t)}$ . Since the ranking of the function value only depends on  $\mathbf{e}^T(X_i - \mathbf{m}^{(t)})/\sigma^{(t)}$ , one may rewrite the update of the mean vector as

$$\mathbf{m}^{(t+1)} = \mathbf{m}^{(t)} + c_m \sigma^{(t)} \sum_{i=1}^{\lambda} w_i \mathcal{N}_{i:\lambda}(0, 1) \cdot \mathbf{e} + c_m \sigma^{(t)} \mu_w^{-\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I} - \mathbf{e}\mathbf{e}^T) , \tag{35}$$

where  $\mathcal{N}_{i:\lambda}(0, 1)$  are the  $i$ -th order statistics from  $\lambda$  population of  $\mathcal{N}(0, 1)$ , and  $\mathcal{N}(\mathbf{0}, \mathbf{I} - \mathbf{e}\mathbf{e}^T)$  is the normally distributed random vector with mean vector  $\mathbf{0}$  and the degenerated covariance matrix  $\mathbf{I} - \mathbf{e}\mathbf{e}^T$ . It indicates that the mean vector moves along the gradient direction with the distribution  $c_m \sigma^{(t)} \sum_{i=1}^{\lambda} w_i \mathcal{N}_{i:\lambda}(0, 1)$ , while it moves randomly in the subspace orthogonal to the gradient with the distribution  $c_m \sigma^{(t)} \mu_w^{-\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I} - \mathbf{e}\mathbf{e}^T)$ .

If the function is spherical, i.e.  $\mathbf{A} \propto \mathbf{I}$ , the mean vector does a symmetric, unbiased random walk on the surface of a hypersphere while the radius of the hypersphere gradually decreases due to the second term on (35). If the function is a

<sup>7</sup> However, the worst case scenario on the discus type function,  $1/(\alpha + (N - 1))$ , describes an empirical observation [36] that the convergence speed of evolution strategy with isotropic distribution does not scale down with  $N$  for  $N \ll \alpha$ .

general convex quadratic function,  $\mathbf{A} \not\propto \mathbf{I}$ , the corresponding random walk on the surface of a hyperellipsoid becomes biased. Then,  $\mathbf{m} - \mathbf{x}^*$  tends to be parallel to the eigenspace corresponding to the smallest eigenvalue  $d_N(\mathbf{A})$ , which means that the quality gain is close to the worst case of  $d_N(\mathbf{A})/\text{Tr}(\mathbf{A})$ . The reason may be explained as follows. The progress in one step is the largest in the short axis direction (parallel to the eigenvector corresponding to the largest eigenvalue of  $\mathbf{A}$ ), and the smallest in the long axis direction (parallel to the eigenvector corresponding to the largest eigenvalue of  $\mathbf{A}$ ). The short axis direction is quickly optimized and the situation gets close to the worst case, where it takes many iterations to escape from. Therefore, we observe the near worst situation in practice. Further theoretical investigation on the distribution of  $\mathbf{e} = \nabla f(\mathbf{m}^{(t)}) / \|\nabla f(\mathbf{m}^{(t)})\|$  should be done in the future work.

#### 4.5. Experiments

To see the effect of the eigenvalue distribution of  $\mathbf{A}$ , we run the experiments. Four quadratic functions are considered: Sphere, Discus, Ellipsoid, Cigar functions of  $N = 10, 100, 1000$  dimensions. The ES with the weights optimal for the infinite dimensional sphere, (8), and the optimal normalized step-size  $\bar{\sigma}^*$  derived for  $c_m \rightarrow \infty$ , (27), times a constant factor is run for  $T = 10000$  iterations. The empirical normalized quality gain is estimated as  $(2/T) \sum_{t=T/2}^{T-1} [f(\mathbf{m}^{(t)}) - f(\mathbf{m}^{(t+1)})] / [f(\mathbf{m}^{(t)})g(\mathbf{m}^{(t)})]$ . The mean vector is initialized randomly by the normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . Eleven independent runs are conducted for each setting. The results are compared with  $\varphi$ , which is supposed to approximate the empirical normalized quality gain for  $c_m \gg 1$  and  $N \gg 1$ . Note that  $\bar{\sigma}^*$  in (27) and  $\varphi$  in (17) depend on  $\mathbf{m}$  through  $\mathbf{e}^T \mathbf{A} \mathbf{e} / \text{Tr}(\mathbf{A})$ . We replace  $\mathbf{e}^T \mathbf{A} \mathbf{e} / \text{Tr}(\mathbf{A})$  with  $d_N(\mathbf{A}) / \text{Tr}(\mathbf{A})$  based on the observation and the above discussion that the mean vector tends to be parallel to the eigenspace corresponding to the smallest eigenvalue of  $\mathbf{A}$ . Figs. 4 and 5 show the median (marker) and the 10%–90% interval (shaded area) of the empirical normalized quality gain for each  $c_m$  and the theoretically derived normalized quality gain formula discussed above. Note that the shaded area is almost invisible, implying that the number of runs and the number of iterations are sufficient to get accurate estimates.

We first focus on the results with  $c_m = 1$  (the default setting). The empirical normalized quality gain gets closer to the normalized quality gain derived for the infinite dimensional quadratic function as  $N$  increases. The approach of the empirical normalized quality gain to the theory is the fastest for the sphere function ( $\mathbf{A} = \mathbf{I}$ ). For convex quadratic functions with the same condition number of  $\alpha = 10^6$ , the speed of the convergence of the normalized quality gain to  $\varphi$  as  $N \rightarrow \infty$  is the fastest for the cigar function, and the slowest for the discus function. This reflects the upper bound derived in Theorem 7 that depends on the ratio  $\text{Tr}(\mathbf{A}^2) / \text{Tr}(\mathbf{A})^2$ , whose value is summarized in Table 1. For the cigar function  $\text{Tr}(\mathbf{A}^2) / \text{Tr}(\mathbf{A})^2$  is close to  $1/(N-1)$ , while for the discus function it is very close to 1 for  $N \ll \alpha$  and we do not observe significant difference between results on different  $N$ .

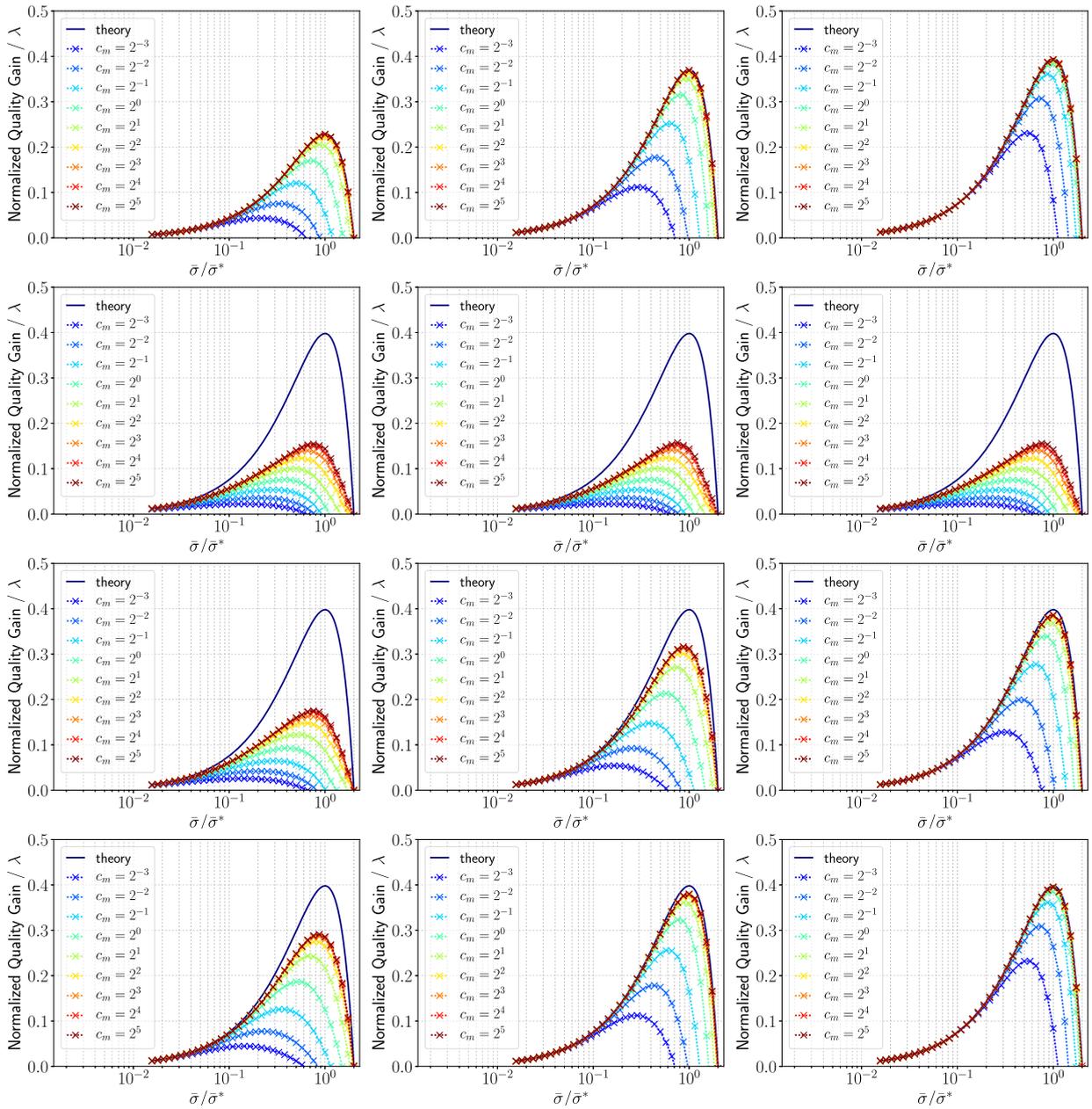
A larger  $c_m$  led to a better empirical normalized quality gain for all cases, i.e., the empirical normalized quality gains became monotonically closer to the theoretical curve.<sup>8</sup> As  $c_m$  becomes greater while the normalized step-size is fixed, the ratio  $\sigma / \|\mathbf{m}\|$  becomes smaller and tends to zero in the limit  $c_m \rightarrow \infty$ . As Corollary 15 implies, the normalized quality gain converges to  $\varphi$  in the limit  $\sigma / \|\mathbf{m}\| \rightarrow 0$ . Therefore, the results reflect the theory. Moreover, the theoretically optimal normalized step-size  $\bar{\sigma}^*$  well approximates the empirically optimal normalized step-size  $\bar{\sigma}$  that maximize the normalized quality gain for all cases when  $c_m \geq 1$ . As  $c_m$  becomes smaller, the empirically optimal normalized step-size  $\bar{\sigma}$  becomes smaller compared to  $\bar{\sigma}^*$ . Note that the difference of the empirical normalized quality gain curves on the sphere function comes only from the randomness of the length of each step  $Z$ . If we replace  $Z$  with  $(\mathbb{E}[\|Z\|] / \|Z\|)Z$  in the algorithm, the selection is independent of  $c_m$  values and is determined by the inner product of the step and the gradient of the objective function at the mean vector. Then, the effect of  $c_m$  goes away.

Comparing Fig. 4 and Fig. 5, the empirical curves are closer to the theoretical curves in Fig. 4. It reflects the fact that the bound between the normalized quality gain and the asymptotic normalized quality gain derived in Theorem 7 typically increases as  $\lambda$  increases. To approximate the theoretical curve, a larger  $c_m$  value is required when  $\lambda$  is greater. The peak of the empirical curves tend to be achieved at a smaller normalized step-size as  $\lambda$  or  $c_m$  becomes greater or smaller, respectively.

## 5. Conclusion

We perform the quality gain analysis of the weighted recombination evolution strategy (ES) on a convex quadratic function. Differently from the previous works, where the limit for the search space dimension  $N$  to infinity is considered, we derive the error bound between the so-called normalized quality gain and its limit expression for the finite dimension. We show that the bound converges to zero when (I)  $N \rightarrow \infty$  as long as the Hessian  $\mathbf{A}$  of the objective function satisfies  $\text{Tr}(\mathbf{A}^2) / \text{Tr}(\mathbf{A})^2 \rightarrow 0$ , or when (II)  $\sigma / \|\mathbf{m}\| \rightarrow 0$ . The limit expression of the normalized quality gain reveals that the optimal recombination weights are independent of the Hessian matrix in the limit (I). Moreover, if the effective variance selection

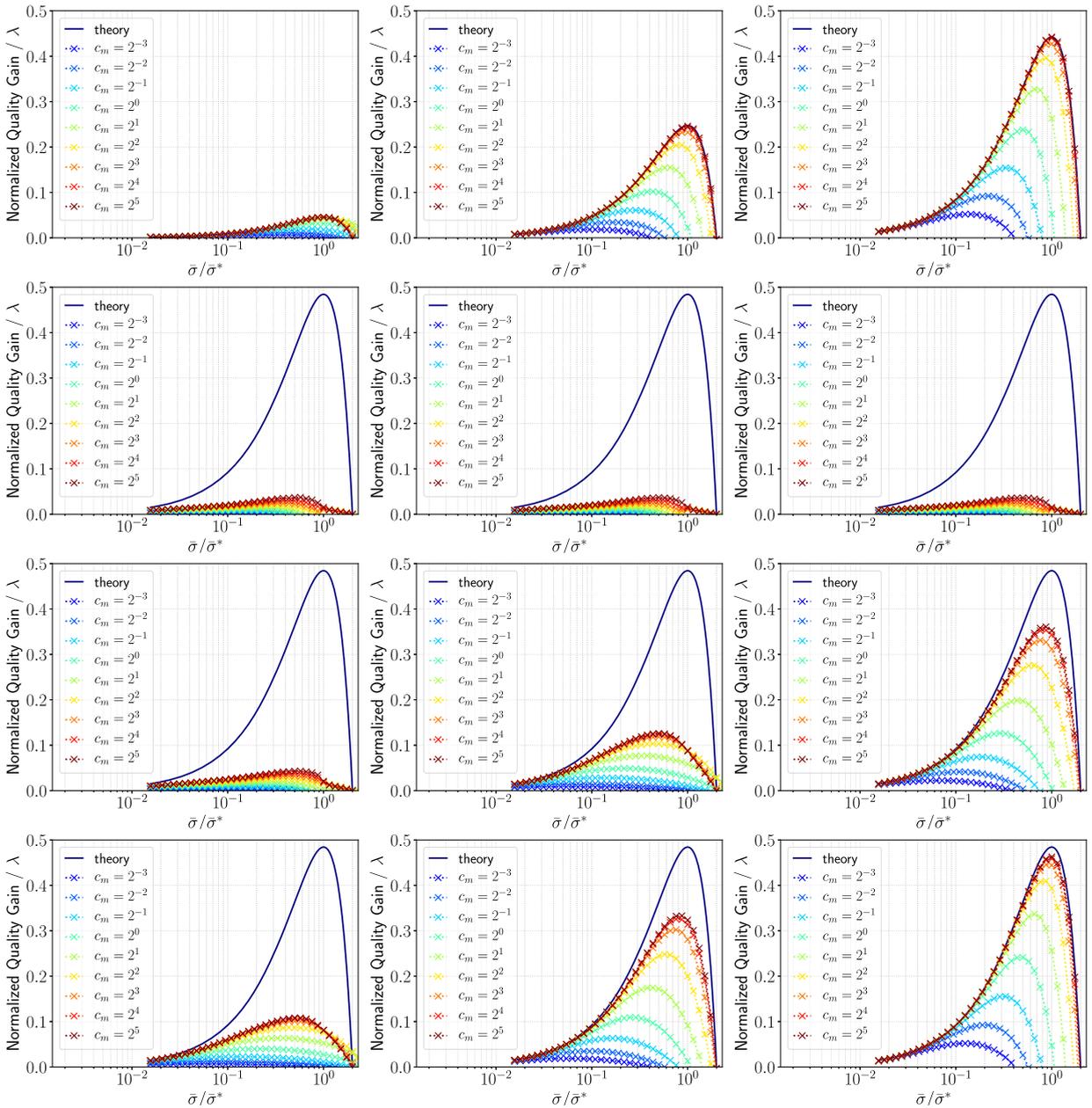
<sup>8</sup> Figure 4 in the previous work [1] shows non-monotonic change of empirical normalized quality gain over  $c_m$ , whereas in Figs. 4 and 5 of this paper shows a monotonic behavior. In Figure 4 in [1]  $\bar{\sigma}^*$  is approximated with (7), whereas in the figures of this paper  $\bar{\sigma}^*$  is computed with (27). The difference between these two quantities is less pronounced as  $N$  increases. The monotonic changes of the graphs are because  $\bar{\sigma}^*$  in (27) approximates the optimal normalized step-size better than (7) on a finite dimensional quadratic function.



**Fig. 4.** Empirical normalized quality gain on four convex quadratic functions, Sphere, Discus, Ellipsoid and Cigar (from top to bottom) of dimension  $N = 10, 100$  and  $1000$  (from left to right). The optimal weights (8) are used and  $\lambda = 10$ .

mass  $\mu_w$  is sufficiently large, the optimal recombination weights for the limit (II) admits the same optimal recombination weights. The optimal normalized step-size for given recombination weights is derived. In the limit (I) the optimal normalized step-size is independent of  $\mathbf{A}$ , while the optimal step-size is proportional to the length of the gradient at the distribution mean. The limit (II) reveals the dependencies of the normalized step-size on  $N$  and  $\mu_w$ .

The quality gain analysis provides a useful insight into the algorithmic behavior, even though it does not take into account the adaptation of the step-size. Knowing the optimal recombination weights ( $w_k^*$ ) directly contributes to the optimal parameter setting. On the contrary, knowing the optimal normalized step-size  $\bar{\sigma}^*$  does not lead to the optimal step-size control. This is because the optimal scale-invariant step-size  $\sigma^*$  in Definition 3 where  $\bar{\sigma}$  is replaced with its optimal value  $\bar{\sigma}^*$  is proportional to  $\|\nabla f(\mathbf{m})\|$ , which is unknown to the algorithm. The optimal step-size, however, is useful to evaluate step-size control mechanisms and to see how close to the optimal situation the step-size control mechanism is. Some theoretical insights into the adaptation mechanism of practical step-size adaptive methods is provided by the approached referred to as “dynamical system” approach by its authors. We refer to [27,28] for the recent development in the dynamical



**Fig. 5.** Empirical normalized quality gain on four convex quadratic functions, Sphere, Discus, Ellipsoid and Cigar (from top to bottom) of dimension  $N = 10, 100$  and  $1000$  (from left to right). The optimal weights (8) are used and  $\lambda = 100$ .

system approach. An important remaining question is: what is the optimal parameter update? Neither the quality gain analysis nor the dynamical system approach will answer this question. The optimal step-size on a quadratic function is revealed in this paper, however, it depends on the norm of the gradient, which is unknown to the real algorithm. A methodology to analyze the optimal update, rather than the optimal parameter value, hopefully including the covariance matrix update is desired in future work.

**Acknowledgements**

The authors thank *Dagstuhl seminar 17191: Theory of Randomized Optimization Heuristics* for providing the opportunity to present and discuss this work. This work is partially supported by JSPS KAKENHI Grant Number 15K16063.

**Appendix A. Normal order statistics**

Here we summarize some important properties of the moments of normal order statistics that are useful to understand the results in the paper.

The first moments of the normal order statistics have the properties:  $\mathbb{E}[\mathcal{N}_{i:\lambda}] \leq \mathbb{E}[\mathcal{N}_{i+1:\lambda}]$ ,  $\mathbb{E}[\mathcal{N}_{i:\lambda}] = -\mathbb{E}[\mathcal{N}_{\lambda+1-i:\lambda}]$ , and  $\sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}] = 0$ . The second (product) moments of the normal order statistics have the following properties:  $\sum_{j=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] = 1$ ,  $\sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}^2] = \sum_{i=1}^{\lambda} \sum_{j=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] = \lambda$ , and  $\mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] = \mathbb{E}[\mathcal{N}_{j:\lambda} \mathcal{N}_{i:\lambda}] = \mathbb{E}[\mathcal{N}_{\lambda+1-i:\lambda} \mathcal{N}_{\lambda+1-j:\lambda}] = \mathbb{E}[\mathcal{N}_{\lambda+1-j:\lambda} \mathcal{N}_{\lambda+1-i:\lambda}]$ .

Here we summarize useful inequalities about order statistics that are all listed in Section 35.1.6 of [37]. The positive dependency inequality tells that the order statistics are non-negatively correlated,  $\text{Cov}(\mathcal{N}_{i:\lambda}, \mathcal{N}_{j:\lambda}) = \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] - \mathbb{E}[\mathcal{N}_{i:\lambda}] \mathbb{E}[\mathcal{N}_{j:\lambda}] \geq 0$ . Together with  $\sum_{j=1}^{\lambda} \text{Cov}(\mathcal{N}_{i:\lambda}, \mathcal{N}_{j:\lambda}) = \sum_{j=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] = 1$ , we have  $0 \leq \text{Cov}(\mathcal{N}_{i:\lambda}, \mathcal{N}_{j:\lambda}) \leq 1$ . It implies  $\mathbb{E}[\mathcal{N}_{i:\lambda}] \mathbb{E}[\mathcal{N}_{j:\lambda}] \leq \mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}] \leq \mathbb{E}[\mathcal{N}_{i:\lambda}] \mathbb{E}[\mathcal{N}_{j:\lambda}] + 1$ .

Another important inequality is David inequality for normal distribution. It tells that  $\Phi^{-1}(i/(\lambda + 1)) \leq \mathbb{E}[\mathcal{N}_{i:\lambda}] \leq \min\{\Phi^{-1}(i/(\lambda + 0.5)), \Phi^{-1}((i - 0.5)/\lambda)\}$ , where  $\Phi$  is the c.d.f. of  $\mathcal{N}(0, 1)$ . It proves an asymptotically tight approximation (Blom’s approximation)  $\mathbb{E}[\mathcal{N}_{i:\lambda}] \approx \Phi^{-1}(\frac{i-\alpha}{\lambda-2\alpha+1})$  with  $\alpha = 0.375$  for  $i \leq \lceil \lambda/2 \rceil$ . The following asymptotic equalities are also used (see Example 8.1.1 in [37])

$$\lim_{\lambda \rightarrow \infty} \frac{\mathbb{E}[\mathcal{N}_{\lambda:\lambda}] - \mathbb{E}[\mathcal{N}_{1:\lambda}]}{2(2 \ln(\lambda))^{\frac{1}{2}}} = 1, \quad \lim_{\lambda \rightarrow \infty} \frac{1}{\lambda} \sum_{i=1}^{\lambda} |\mathbb{E}[\mathcal{N}_{i:\lambda}]| = \frac{2^{\frac{1}{2}}}{\pi^{\frac{1}{2}}}, \quad \lim_{\lambda \rightarrow \infty} \frac{1}{\lambda} \sum_{i=1}^{\lambda} \mathbb{E}[\mathcal{N}_{i:\lambda}]^2 = 1. \tag{A.1}$$

Let  $\mathbf{n}_{(\lambda)}$  be the  $\lambda$  dimensional column vector whose  $i$ -th component is  $\mathbb{E}[\mathcal{N}_{i:\lambda}]$  and  $\mathbf{N}_{(\lambda)}$  be the  $\lambda$  dimensional symmetric matrix whose  $(i, j)$ -th element is  $\mathbb{E}[\mathcal{N}_{i:\lambda} \mathcal{N}_{j:\lambda}]$ . The covariance matrix  $\mathbf{N}_{(\lambda)} - \mathbf{n}_{(\lambda)} \mathbf{n}_{(\lambda)}^T$  is by definition nonnegative definite. It implies the eigenvalues of  $\mathbf{N}_{(\lambda)}$  are all nonnegative. Moreover, from the above mentioned fact derives that the sum of the eigenvalues is  $\text{Tr}(\mathbf{N}_{(\lambda)}) = \sum_{i=1}^{\lambda} \sum_{j=1}^{\lambda} \text{Cov}(\mathcal{N}_{i:\lambda}, \mathcal{N}_{j:\lambda}) = \lambda$ . Furthermore, the third asymptotic relation of (A.1) reads  $\lim_{\lambda \rightarrow \infty} \text{Tr}(\mathbf{n}_{(\lambda)} \mathbf{n}_{(\lambda)}^T) / \lambda = \lim_{\lambda \rightarrow \infty} \|\mathbf{n}_{(\lambda)}\|^2 / \lambda = 1$ . It implies, for any  $\mathbf{x} \in \mathbb{R}^{\lambda} \setminus \{\mathbf{0}\}$ , we have

$$\lim_{\lambda \rightarrow \infty} \frac{\mathbf{x}^T \mathbf{N}_{(\lambda)} \mathbf{x}}{\lambda \|\mathbf{x}\|^2} = \lim_{\lambda \rightarrow \infty} \frac{(\mathbf{x}^T \mathbf{n}_{(\lambda)})^2}{\lambda \|\mathbf{x}\|^2} = \lim_{\lambda \rightarrow \infty} \frac{(\mathbf{x}^T \mathbf{n}_{(\lambda)})^2}{\|\mathbf{x}\|^2 \|\mathbf{n}_{(\lambda)}\|^2}. \tag{A.2}$$

**Appendix B. Proofs and derivations**

*B.1. Proof of Proposition 2*

**Proof.** Let  $\Delta = c_m \sum_{i=1}^{\lambda} W(i; (\mathbf{m} + \sigma Z_k)_{k=1}^{\lambda}) Z_i$ , where  $(Z_i)_{i=1}^{\lambda}$  are independent and  $N$ -variate standard normally distributed random vectors. Then,

$$\begin{aligned} \phi(\mathbf{m}, \sigma) &= 1 - \mathbb{E}[f^*(\mathbf{m} + \sigma \Delta)] / f^*(\mathbf{m}) \\ &= 1 - \mathbb{E}[f(\mathbf{m} + \sigma \Delta - \mathbf{x}^*)] / f(\mathbf{m} - \mathbf{x}^*) \\ &= 1 - \alpha^{-n} \mathbb{E}[f(\alpha \cdot (\mathbf{m} + \sigma \Delta - \mathbf{x}^*))] / \alpha^{-n} f(\alpha \cdot (\mathbf{m} - \mathbf{x}^*)) \\ &= 1 - \mathbb{E}[f(\alpha \cdot (\mathbf{m} + \sigma \Delta - \mathbf{x}^*))] / f(\alpha \cdot (\mathbf{m} - \mathbf{x}^*)) \\ &= 1 - \mathbb{E}[f^*(\mathbf{x}^* + \alpha \cdot (\mathbf{m} - \mathbf{x}^*) + \alpha \sigma \Delta)] / f^*(\mathbf{x}^* + \alpha \cdot (\mathbf{m} - \mathbf{x}^*)) \\ &= \phi(\mathbf{x}^* + \alpha(\mathbf{m} - \mathbf{x}^*), \alpha \sigma). \end{aligned}$$

Note that  $\phi(\mathbf{x}^* + (\mathbf{m} - \mathbf{x}^*), \sigma) = \phi(\mathbf{m}, \sigma)$ . That is, the quality gain is scale invariant around  $(\mathbf{x}^*, 0)$ . Moreover, the above equality implies that  $\text{argmax}_{\sigma} \phi(\mathbf{x}^* + (\mathbf{m} - \mathbf{x}^*), \sigma) = \text{argmax}_{\sigma} \phi(\mathbf{x}^* + \alpha(\mathbf{m} - \mathbf{x}^*), \alpha \sigma)$ , i.e., the optimal step-size at  $\mathbf{x}^* + \alpha(\mathbf{m} - \mathbf{x}^*)$  is  $\alpha$  times greater than the optimal step-size at  $\mathbf{x}^* + (\mathbf{m} - \mathbf{x}^*)$ . Therefore, the optimal step-size as a function of  $\mathbf{m} - \mathbf{x}^*$  is homogeneous of degree 1, i.e.,  $\sigma^*(\alpha \cdot (\mathbf{m} - \mathbf{x}^*)) = \alpha \sigma^*(\mathbf{m} - \mathbf{x}^*)$ .  $\square$

*B.2. Proof of Lemma 5*

**Proof.** Since  $(X_k)_{k=1}^{\lambda}$  are independent and normally distributed, the conditional probability of  $\mathbb{1}_{f(X_k) < f(X_i)} = 1$  given  $X_i$  for any  $k \neq i$  is  $F_f(f(X_i))$ . Then, the probability of  $\sum_{k=1}^{\lambda} \mathbb{1}_{f(X_k) \leq f(X_i)}$  being  $a$  for  $a \in \llbracket 1, \lambda \rrbracket$  is given by  $P_b(a - 1; \lambda - 1, p)$  with  $p = F_f(f(X_i))$ . Then, for any  $\alpha \geq 0$ ,

$$\mathbb{E}_i[W(i; (X_k)_{k=1}^{\lambda})^{\alpha}] = \sum_{k=1}^{\lambda} w_k^{\alpha} P_b(k - 1; \lambda - 1, p).$$

Similarly, the joint distribution of  $\sum_{k=1}^{\lambda} \mathbb{1}_{f(X_k) \leq f(X_i)}$  and  $\sum_{k=1}^{\lambda} \mathbb{1}_{f(X_k) \leq f(X_j)}$  is derived. Due to the symmetry between  $i$  and  $j$ , we can assume w.l.o.g. that  $f(X_i) \leq f(X_j)$ . Then, the joint probability of  $\sum_{k=1}^{\lambda} \mathbb{1}_{f(X_k) \leq f(X_i)} = a$  and  $\sum_{k=1}^{\lambda} \mathbb{1}_{f(X_k) \leq f(X_j)} = b$  for  $a, b \in \llbracket 1, \lambda \rrbracket$  is given by  $P_t(a-1, b-a-1; \lambda-2, p, q-p)$  with  $p = F_f(f(X_i))$  and  $q = F_f(f(X_j))$  if  $a < b$ , and zero otherwise. Then,

$$\mathbb{E}_{i,j}[W(i; (X_k)_{k=1}^{\lambda})W(j; (X_k)_{k=1}^{\lambda})] = \sum_{m=1}^{\lambda-1} \sum_{l=m+1}^{\lambda} w_m w_l P_t(m-1, l-m-1; \lambda-2, p, q-p) .$$

This ends the proof.  $\square$

### B.3. Proof of Lemma 6

**Proof.** The derivative of  $u_1$  is  $\sum_{k=1}^{\lambda} w_k \binom{\lambda-1}{k-1} \frac{d}{dp} [p^{k-1}(1-p)^{\lambda-k}]$ , where

$$\frac{d}{dp} [p^{k-1}(1-p)^{\lambda-k}] = (k-1)p^{k-2}(1-p)^{\lambda-k} - (\lambda-k)p^{k-1}(1-p)^{\lambda-k-1} .$$

Substituting the derivatives and rearranging the terms, we obtain

$$\frac{du_1(p)}{dp} = (\lambda-1) \sum_{k=1}^{\lambda-1} (w_{k+1} - w_k) \binom{\lambda-2}{k-1} p^{k-1}(1-p)^{\lambda-k-1} .$$

The Lipschitz constant  $L_1$  is the supremum of the absolute value of the derivative derived above. It completes the proof for the  $\ell_1$ -Lipschitz continuity of  $u_1$  and its Lipschitz constant. Since  $u_2$  is equivalent to  $u_1$  if  $w_i$  are replaced with  $w_i^2$  in the definition of  $u_1$ , we have the  $\ell_1$ -Lipschitz continuity of  $u_2$  and its Lipschitz constant by replacing  $w_i$  with  $w_i^2$  in the above argument.

The partial derivative of  $u_3$  with respect to  $p$  is

$$\sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \binom{\lambda-2}{l-2} \binom{l-2}{k-1} \frac{\partial}{\partial p} [\min(p, q)^{k-1} |q-p|^{l-k-1} (1-\min(p, q))^{\lambda-l}] ,$$

where

$$\begin{aligned} & \frac{\partial}{\partial p} \min(p, q)^{k-1} |q-p|^{l-k-1} (1-\min(p, q))^{\lambda-l} \\ &= \begin{cases} [(k-1)(q-p) - (l-k-1)p] p^{k-2} (q-p)^{l-k-2} (1-q)^{\lambda-l} & (p < q) \\ [(l-k-1)(1-p) - (\lambda-l)(p-q)] q^{k-1} (p-q)^{l-k-2} (1-p)^{\lambda-l-1} & (p > q) \end{cases} . \end{aligned}$$

Substituting the derivatives and rearranging the terms, we obtain

$$\begin{aligned} & \frac{1}{\lambda-2} \frac{\partial u_3(p, q)}{\partial p} \\ &= \begin{cases} \sum_{k=1}^{\lambda-2} \sum_{l=k+2}^{\lambda} w_l (w_{k+1} - w_k) \binom{\lambda-3}{l-3} \binom{l-3}{k-1} p^{k-1} (q-p)^{l-k-2} (1-q)^{\lambda-l} & (p < q) \\ \sum_{k=1}^{\lambda-2} \sum_{l=k+2}^{\lambda} w_k (w_l - w_{l-1}) \binom{\lambda-3}{l-3} \binom{l-3}{k-1} q^{k-1} (p-q)^{l-k-2} (1-p)^{\lambda-l} & (p > q) \end{cases} \end{aligned}$$

Since  $u_3$  is differentiable with respect to  $p$  almost everywhere in  $(0, 1)$ , it is Lipschitz continuous with respect to  $p$ . Its Lipschitz constant is  $\sup_{q \in (0,1)} \sup_{p \in (0,q) \cup (q,1)} \left| \frac{\partial u_3(p,q)}{\partial p} \right|$ . Due to the symmetry,  $u_3(p, q)$  is  $\ell_1$ -Lipschitz continuous on  $[0, 1]^2$  with the Lipschitz constant  $L_3 = \sup_{q \in (0,1)} \sup_{p \in (0,q) \cup (q,1)} \left| \frac{\partial u_3(p,q)}{\partial p} \right|$ . This completes the proof.  $\square$

### B.4. Upper bounds of Lipschitz constants

For a general weight scheme, we have the following trivial upper bounds for the Lipschitz constants derived in Lemma 6,

$$L_1 \leq (\lambda-1) \max_{k \in \llbracket 1, \lambda-1 \rrbracket} |w_{k+1} - w_k| , \tag{B.1}$$

$$L_2 \leq (\lambda-1) \max_{k \in \llbracket 1, \lambda-1 \rrbracket} |w_{k+1}^2 - w_k^2| , \tag{B.2}$$

$$L_3 \leq (\lambda-2) \max_{k \in \llbracket 1, \lambda \rrbracket} \max_{l \in \llbracket 1, k-2 \rrbracket \cup \llbracket k+1, \lambda-1 \rrbracket} |w_k| \cdot |w_{l+1} - w_l| . \tag{B.3}$$

These upper bounds are straight-forward from the facts  $\sum_{k=1}^{\lambda-1} P_b(k-1; \lambda-2, p) = 1$  and  $\sum_{k=1}^{\lambda-2} \sum_{l=k+2}^{\lambda} P_t(k-1, l-k-2; \lambda-3, \min(p, q), |q-p|) = 1$ .

For the truncation weights with  $3 \leq \mu \leq \lambda-2$ , we can obtain better bounds. The bounds of the factorial of  $n \geq 1$  known by Robbins [38], namely,

$$(2\pi n)^{\frac{1}{2}} \left(\frac{n}{e}\right)^n \exp\left(\frac{1}{12n+1}\right) < n! < (2\pi n)^{\frac{1}{2}} \left(\frac{n}{e}\right)^n \exp\left(\frac{1}{12n}\right)$$

gives us an upper bound of  $\binom{n}{k}$  for  $0 < k < n$

$$\binom{n}{k} < \left(\frac{n}{2\pi k(n-k)}\right)^{\frac{1}{2}} \left(\frac{n}{k}\right)^k \left(\frac{n}{n-k}\right)^{n-k} . \tag{B.4}$$

Here we used  $\exp\left(\frac{1}{12n} - \frac{1}{12k+1} - \frac{1}{12(n-k)+1}\right) < 1$ . On the other hand, we have for  $0 < k < n$

$$\sup_{0 \leq p \leq 1} p^k(1-p)^{n-k} = \left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} . \tag{B.5}$$

Since  $w_{k+1} - w_k = -1/\mu$  for  $k = \mu$  and  $w_{k+1} - w_k = 0$  for  $k \neq \mu$ , we have for  $3 \leq \mu \leq \lambda-2$ ,

$$\begin{aligned} L_1 &= \sup_{0 < p < 1} \left| (\lambda-1) \frac{1}{\mu} \binom{\lambda-2}{\mu-1} p^{\mu-1} (1-p)^{\lambda-\mu-1} \right| \\ &= \frac{\lambda-1}{\mu} \binom{\lambda-2}{\mu-1} \left(\frac{\mu-1}{\lambda-2}\right)^{\mu-1} \left(\frac{\lambda-\mu-1}{\lambda-2}\right)^{\lambda-\mu-1} \\ &\leq \frac{\lambda-1}{\mu} \left(\frac{\lambda-2}{2\pi(\mu-1)(\lambda-\mu-1)}\right)^{\frac{1}{2}} . \end{aligned}$$

Analogously, since  $w_{k+1}^2 - w_k^2 = -1/\mu^2$  for  $k = \mu$  and  $w_{k+1}^2 - w_k^2 = 0$  for  $k \neq \mu$ , we obtain the bound of  $L_2$ :  $L_2 < [(\lambda-1)/\mu^2] \cdot [(\lambda-2)/(2\pi(\mu-1)(\lambda-\mu-1))]^{1/2}$ .

Moreover, since  $w_k(w_l - w_{l-1}) = -1/\mu^2$  for  $l = \mu+1$  and  $w_k(w_l - w_{l-1}) = 0$  otherwise, we have

$$\begin{aligned} L_3 &= \sup_{q \in (0,1)} \sup_{p \in (q,1)} \sum_{k=1}^{\mu-1} \frac{\lambda-2}{\mu^2} \binom{\lambda-3}{\mu-2} \binom{\mu-2}{k-1} q^{k-1} (p-q)^{\mu-k-1} (1-p)^{\lambda-\mu-1} \\ &= \sup_{p \in (0,1)} \frac{(\lambda-2)}{\mu^2} \binom{\lambda-3}{\mu-2} (1-p)^{\lambda-\mu-1} \sup_{q \in (0,p)} \sum_{k=1}^{\mu-1} \binom{\mu-2}{k-1} q^{k-1} (p-q)^{\mu-k-1} \\ &= \frac{(\lambda-2)}{\mu^2} \binom{\lambda-3}{\mu-2} \sup_{p \in (0,1)} (1-p)^{\lambda-\mu-1} p^{\mu-2} \\ &= \frac{(\lambda-2)}{\mu^2} \binom{\lambda-3}{\mu-2} \left(\frac{\mu-2}{\lambda-3}\right)^{\mu-2} \left(\frac{\lambda-\mu-1}{\lambda-3}\right)^{\lambda-\mu-1} \\ &\leq \frac{(\lambda-2)}{\mu^2} \left(\frac{\lambda-3}{2\pi(\mu-2)(\lambda-\mu-1)}\right)^{\frac{1}{2}} . \end{aligned}$$

Here we used (B.4), (B.5), and the binomial relation  $\sum_{k=1}^{\mu-1} \binom{\mu-2}{k-1} q^{k-1} (p-q)^{\mu-k-1} = p^{\mu-2}$ .

**B.5. Proof of Lemma 8**

**Proof.** If  $\alpha = 1$ , then  $G(\alpha) = 1$ , and the inequality is trivial. Hence, we assume  $\alpha < 1$  in the following.

Remember that  $H_N = Z_e + h(Z)$ . The absolute difference between  $F_N(t)$  and  $\Phi(t)$  is rewritten as follows

$$\begin{aligned} |F_N(t) - \Phi(t)| &= |\Pr[H_N \leq t] - \Pr[Z_e \leq t]| \\ &= |\Pr[Z_e + h(Z) \leq t] - \Pr[Z_e \leq t]| \\ &= \Pr[h(Z) \geq 0 \text{ and } t - h(Z) \leq Z_e \leq t] + \Pr[h(Z) \leq 0 \text{ and } t \leq Z_e \leq t - h(Z)] . \end{aligned}$$

With an arbitrary  $\epsilon_+ > 0$ , the first term on the RMS is upper bounded as

$$\begin{aligned} & \Pr[h(Z) \geq 0 \text{ and } t - h(Z) \leq Z_{\mathbf{e}} \leq t] \\ & \leq \Pr[h(Z) \geq \epsilon_+] + \Pr[h(Z) < \epsilon_+ \text{ and } t - h(Z) \leq Z_{\mathbf{e}} \leq t] \\ & \leq \Pr[h(Z) \geq \epsilon_+] + \Pr[h(Z) < \epsilon_+ \text{ and } t - \epsilon_+ \leq Z_{\mathbf{e}} \leq t] \\ & \leq \Pr[h(Z) \geq \epsilon_+] + \Pr[t - \epsilon_+ \leq Z_{\mathbf{e}} \leq t] \\ & \leq \Pr[h(Z) \geq \epsilon_+] + (2\pi)^{-\frac{1}{2}} \epsilon_+ . \end{aligned}$$

For the last inequality, we used that the density of the one-dimensional standard normal distribution is at most  $(2\pi)^{-\frac{1}{2}}$  and  $Z_{\mathbf{e}}$  is of the one-dimensional standard normal distribution. Analogously, we have for any  $\epsilon_- > 0$

$$\begin{aligned} \Pr[h(Z) \leq 0 \text{ and } t \leq Z_{\mathbf{e}} \leq t - h(Z)] & \leq \Pr[h(Z) \leq -\epsilon_-] + \Pr[t \leq Z_{\mathbf{e}} \leq t + \epsilon_-] \\ & \leq \Pr[h(Z) \leq -\epsilon_-] + (2\pi)^{-\frac{1}{2}} \epsilon_- . \end{aligned}$$

Let  $\tilde{h}(Z) = 2(c_m/\bar{\sigma})h(Z) = Z^T \mathbf{A} Z - 1$ ,  $\tilde{\epsilon}_+ = 2(c_m/\bar{\sigma})\epsilon_+$  and  $\tilde{\epsilon}_- = 2(c_m/\bar{\sigma})\epsilon_-$ . Then,  $\Pr[h(Z) \geq \epsilon_+] = \Pr[\tilde{h}(Z) \geq \tilde{\epsilon}_+]$  and  $\Pr[h(Z) \leq -\epsilon_-] = \Pr[\tilde{h}(Z) \leq -\tilde{\epsilon}_-]$ . From Lemma 1 in [32] knows that for any  $x \geq 0$

$$\begin{aligned} \Pr[\tilde{h}(Z) \geq 2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} x^{\frac{1}{2}} + 2d_1(\mathbf{A})x] & \leq \exp(-x) , \\ \Pr[\tilde{h}(Z) \leq -2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} x^{\frac{1}{2}}] & \leq \exp(-x) . \end{aligned} \tag{B.6}$$

Let  $x = \ln(1/\alpha)$  and let  $\epsilon_+$  and  $\epsilon_-$  such that

$$\begin{aligned} \tilde{\epsilon}_+ & = 2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} x^{\frac{1}{2}} + 2d_1(\mathbf{A})x = 2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} ((\ln(1/\alpha))^{\frac{1}{2}} + (d_1(\mathbf{A})/\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}) \ln(1/\alpha)) \\ \tilde{\epsilon}_- & = 2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} x^{\frac{1}{2}} = 2 \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} (\ln(1/\alpha))^{\frac{1}{2}} . \end{aligned}$$

Then, from (B.6) derives that

$$\begin{aligned} & \Pr[h(Z) \geq \epsilon_+] + (2\pi)^{-\frac{1}{2}} \epsilon_+ \\ & = \Pr[\tilde{h}(Z) \geq \tilde{\epsilon}_+] + (\bar{\sigma}\tilde{\epsilon}_+)/ (2(2\pi)^{\frac{1}{2}} c_m) \\ & \leq \alpha + (\bar{\sigma}\tilde{\epsilon}_+)/ (2(2\pi)^{\frac{1}{2}} c_m) \\ & = \alpha + (2\pi)^{-\frac{1}{2}} (\bar{\sigma}/c_m) \text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} ((\ln(1/\alpha))^{\frac{1}{2}} + (d_1(\mathbf{A})/\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}) \ln(1/\alpha)) \\ & = \alpha (1 + (2\pi)^{-\frac{1}{2}} (\ln(1/\alpha))^{\frac{1}{2}} + (2\pi)^{-\frac{1}{2}} (d_1(\mathbf{A})/\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}) \ln(1/\alpha)) . \end{aligned}$$

Similarly, we have  $\Pr[h(Z) \leq -\epsilon_-] + (2\pi)^{-\frac{1}{2}} \epsilon_- \leq \alpha (1 + (2\pi)^{-\frac{1}{2}} (\ln(1/\alpha))^{\frac{1}{2}})$ . Altogether, we obtain

$$|F_N(t) - \Phi(t)| \leq \alpha (2 + (2/\pi)^{\frac{1}{2}} (\ln(1/\alpha))^{\frac{1}{2}} + (2\pi)^{-\frac{1}{2}} (d_1(\mathbf{A})/\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}}) \ln(1/\alpha)) .$$

Since the RHS of the above inequality is independent of  $t$ , taking the supremum of both sides over  $t \in \mathbb{R}$ , we obtain the desired inequality.  $\square$

### B.6. Proof of Lemma 9

**Proof.** First, note that  $F_f(f(\mathbf{m} + \sigma Z)) = F_N(H_N)$  and  $H_N = Z_{\mathbf{e}} + h(Z)$ . Using Lemma 6, we have

$$\begin{aligned} |\mathbb{E}[u_1(F_f(f(X)))Z_{\mathbf{e}}] - \mathbb{E}[u_1(\Phi(Z_{\mathbf{e}}))Z_{\mathbf{e}}]| & = |\mathbb{E}[u_1(F_N(H_N))Z_{\mathbf{e}}] - \mathbb{E}[u_1(\Phi(Z_{\mathbf{e}}))Z_{\mathbf{e}}]| \\ & \leq \mathbb{E}[|u_1(F_N(H_N)) - u_1(\Phi(Z_{\mathbf{e}}))| \cdot |Z_{\mathbf{e}}|] \leq L_1 \mathbb{E}[|F_N(H_N) - \Phi(Z_{\mathbf{e}})| \cdot |Z_{\mathbf{e}}|] . \end{aligned}$$

Noting that  $\Phi$  is Lipschitz continuous with the Lipschitz constant  $(2\pi)^{-\frac{1}{2}}$ , we have  $|\Phi(H_N) - \Phi(Z_{\mathbf{e}})| \leq (2\pi)^{-\frac{1}{2}} |H_N - Z_{\mathbf{e}}| = (2\pi)^{-\frac{1}{2}} |h(Z)|$ . On the other hand, Lemma 8 says that  $|F_N(H_N) - \Phi(H_N)| \leq G(\alpha)$ . From these inequalities we obtain

$$|F_N(H_N) - \Phi(Z_{\mathbf{e}})| = |F_N(H_N) - \Phi(H_N) + \Phi(H_N) - \Phi(Z_{\mathbf{e}})| \leq G(\alpha) + (2\pi)^{-\frac{1}{2}} |h(Z)| . \tag{B.7}$$

Using the inequality (B.7) and the Schwarz inequality and the identities  $\mathbb{E}[|Z_{\mathbf{e}}|] = (2/\pi)^{\frac{1}{2}}$ ,  $\mathbb{E}[Z_{\mathbf{e}}^2] = 1$ , and

$$\mathbb{E}[|h(Z)|^2] = \left(\frac{1}{2} \frac{\bar{\sigma}}{c_m}\right)^2 \mathbb{E}[(Z^T \mathbf{A} Z - 1)^2] = \left(\frac{1}{2} \frac{\bar{\sigma}}{c_m}\right)^2 (2 \text{Tr}(\mathbf{A}^2)) = \frac{\alpha^2}{2} , \tag{B.8}$$

we have

$$\begin{aligned} \mathbb{E}[|F_N(H_N) - \Phi(Z_e)| \cdot |Z_e|] &\leq G(\alpha)\mathbb{E}[|Z_e|] + (2\pi)^{-\frac{1}{2}}\mathbb{E}[|h(Z)| \cdot |Z_e|] \\ &\leq G(\alpha)\mathbb{E}[|Z_e|] + (2\pi)^{-\frac{1}{2}}\mathbb{E}[h(Z)^2]^{\frac{1}{2}}\mathbb{E}[Z_e^2]^{\frac{1}{2}} \\ &= (2/\pi)^{\frac{1}{2}}G(\alpha) + (2\pi)^{-\frac{1}{2}}\mathbb{E}[h(Z)^2]^{\frac{1}{2}} \\ &= (2/\pi)^{\frac{1}{2}}G(\alpha) + (4\pi)^{-\frac{1}{2}}\alpha . \end{aligned}$$

Altogether, we obtain the inequality stated in the lemma. This completes the proof.  $\square$

*B.7. Proof of Lemma 10*

**Proof.** Analogously to the proof of Lemma 9, we have

$$\begin{aligned} &|\mathbb{E}[u_2(F_f(f(X)))(Z^T\mathbf{AZ} - 1)] - \mathbb{E}[u_2(\Phi(Z_e))(Z^T\mathbf{AZ} - 1)]| \\ &\leq L_2\mathbb{E}[|F_N(H_N) - \Phi(Z_e)| \cdot |Z^T\mathbf{AZ} - 1|] \\ &\leq L_2\mathbb{E}[(G(\alpha) + (2\pi)^{-\frac{1}{2}}|h(Z)|) |Z^T\mathbf{AZ} - 1|] \\ &= L_2(G(\alpha)\mathbb{E}[|Z^T\mathbf{AZ} - 1|] + (2\pi)^{-\frac{1}{2}}\mathbb{E}[|h(Z)| \cdot |Z^T\mathbf{AZ} - 1|]) . \end{aligned}$$

Applying the inequalities  $\mathbb{E}[|Z^T\mathbf{AZ} - 1|] \leq \mathbb{E}[(Z^T\mathbf{AZ} - 1)^2]^{\frac{1}{2}} = (2\text{Tr}(\mathbf{A}^2))^{\frac{1}{2}}$  and

$$\mathbb{E}[|h(Z)| \cdot |Z^T\mathbf{AZ} - 1|] = \frac{1}{2}(\bar{\sigma}/c_m)\mathbb{E}[(Z^T\mathbf{AZ} - 1)^2] = (\bar{\sigma}/c_m)\text{Tr}(\mathbf{A}^2) = \alpha\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} ,$$

we obtain the inequality stated in the lemma. This completes the proof.  $\square$

*B.8. Proof of Lemma 11*

**Proof.** Using Lemma 6, we have

$$\begin{aligned} &|\mathbb{E}[u_3(F_f(f(X)), F_f(f(\tilde{X})))Z^T\mathbf{A}\tilde{Z}] - \mathbb{E}[u_3(\Phi(Z_e), \Phi(\tilde{Z}_e))Z^T\mathbf{A}\tilde{Z}]| \\ &\leq \mathbb{E}\left[|u_3(F_N(H_N), F_N(\tilde{H}_N)) - u_3(\Phi(Z_e), \Phi(\tilde{Z}_e))| \cdot |Z^T\mathbf{A}\tilde{Z}|\right] \\ &\leq L_3\mathbb{E}\left[ (|F_N(H_N) - \Phi(Z_e)| + |F_N(\tilde{H}_N) - \Phi(\tilde{Z}_e)|) \cdot |Z^T\mathbf{A}\tilde{Z}| \right] . \end{aligned}$$

Then, using the equality  $\mathbb{E}[|Z^T\mathbf{A}\tilde{Z}| | Z] = (2/\pi)^{\frac{1}{2}}\|\mathbf{AZ}\|$  (since  $|Z^T\mathbf{A}\tilde{Z}|$  given  $Z$  is half-normally distributed), the symmetry of  $Z$  and  $\tilde{Z}$ , the Schwarz inequality, and the inequality (B.7), we have

$$\begin{aligned} &\mathbb{E}\left[ (|F_N(H_N) - \Phi(Z_e)| + |F_N(\tilde{H}_N) - \Phi(\tilde{Z}_e)|) \cdot |Z^T\mathbf{A}\tilde{Z}| \right] \\ &\leq 2G(\alpha)\mathbb{E}[|Z^T\mathbf{A}\tilde{Z}|] + 2(2\pi)^{-\frac{1}{2}}\mathbb{E}[|h(Z)| \cdot |Z^T\mathbf{A}\tilde{Z}|] . \end{aligned}$$

On one hand, we have

$$\begin{aligned} \mathbb{E}[|Z^T\mathbf{A}\tilde{Z}|] &= \mathbb{E}[\mathbb{E}[|Z^T\mathbf{A}\tilde{Z}| | Z]] = (2/\pi)^{\frac{1}{2}}\mathbb{E}[\|\mathbf{AZ}\|] \\ &\leq (2/\pi)^{\frac{1}{2}}\mathbb{E}[\|\mathbf{AZ}\|^2]^{\frac{1}{2}} = (2/\pi)^{\frac{1}{2}}\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} , \end{aligned}$$

where we used  $\mathbb{E}[\|\mathbf{AZ}\|^2] = \mathbb{E}[\text{Tr}(\mathbf{AZ}Z^T\mathbf{A})] = \text{Tr}(\mathbf{A}\mathbb{E}[ZZ^T]\mathbf{A}) = \text{Tr}(\mathbf{A}^2)$ . On the other hand, we have

$$\begin{aligned} \mathbb{E}[|h(Z)| \cdot |Z^T\mathbf{A}\tilde{Z}|] &= \mathbb{E}[|h(Z)| \mathbb{E}[|Z^T\mathbf{A}\tilde{Z}| | Z]] = (2/\pi)^{\frac{1}{2}}\mathbb{E}[|h(Z)| \cdot \|\mathbf{AZ}\|] \\ &\leq (2/\pi)^{\frac{1}{2}}\mathbb{E}[|h(Z)|^2]^{\frac{1}{2}}\mathbb{E}[\|\mathbf{AZ}\|^2]^{\frac{1}{2}} = \pi^{-\frac{1}{2}}\alpha\text{Tr}(\mathbf{A}^2)^{\frac{1}{2}} , \end{aligned}$$

where we used  $\mathbb{E}[|h(Z)|^2] = \alpha^2/2$  derived in (B.8). Altogether, we obtain the inequality stated in the lemma. This completes the proof.  $\square$

### B.9. Proof of Lemma 12

**Proof.** Let  $p$  be the probability density function of the one-dimensional standard normal distribution and  $p_{i:\lambda}$  be the probability density function of  $\mathcal{N}_{i:\lambda}$  and  $p_{i,j:\lambda}$  be the joint probability density function of  $\mathcal{N}_{i:\lambda}$  and  $\mathcal{N}_{j:\lambda}$ . It is well known that  $p_{i:\lambda}(x) = \lambda \binom{\lambda-1}{i-1} \Phi(x)^{i-1} (1 - \Phi(x))^{\lambda-i} p(x)$  and  $p_{i,j:\lambda}(x, y) = \lambda(\lambda-1) \binom{\lambda-2}{j-2} \binom{j-1}{i-1} \Phi(x)^{i-1} (\Phi(y) - \Phi(x))^{(j-i-1)} (1 - \Phi(x))^{\lambda-j} p(x)p(y)$  for  $i < j$  and  $x < y$ , and  $p_{i,j:\lambda}(x, y) = 0$  for  $i < j$  and  $x \geq y$ . Note also that  $p_{i,j:\lambda}(x, y) = p_{j,i:\lambda}(y, x)$ .

The functions  $u_1$  and  $u_2$  are then written using these p.d.f.s of the normal order statistics as  $\lambda u_1(\Phi(x))p(x) = \sum_{k=1}^{\lambda} w_k p_{k:\lambda}(x)$  and  $\lambda u_2(\Phi(x))p(x) = \sum_{k=1}^{\lambda} w_k^2 p_{k:\lambda}(x)$ . From these identities, we obtain (19) and (20). The identity (21) is derived by using  $\lambda \mathbb{E}[u_2(\Phi(\mathbf{Z}_e))(Z_e^2 - 1)] = \sum_{i=1}^{\lambda} w_i^2 (\mathbb{E}[\mathcal{N}_{i:\lambda}^2] - 1)$  and  $\mathbb{E}[u_2(\Phi(\mathbf{Z}_e))(Z^T \mathbf{A} \mathbf{Z} - 1)] = \mathbb{E}[u_2(\Phi(\mathbf{Z}_e))(Z_e^2 - 1)] \mathbf{e}^T \mathbf{A} \mathbf{e}$ , where the last equality is proved by using the expression  $Z^T \mathbf{A} \mathbf{Z} = Z_e^2 \mathbf{e}^T \mathbf{A} \mathbf{e} + Z_e \mathbf{e}^T \mathbf{A} \mathbf{Z}_{\perp} + Z_{\perp}^T \mathbf{A} \mathbf{Z}_{\perp}$ , the mutual independence between  $Z_e$  and  $Z_{\perp}$ , and  $\mathbb{E}[Z_{\perp}] = 0$  and  $\mathbb{E}[Z_{\perp}^T \mathbf{A} \mathbf{Z}_{\perp}] = 1 - \mathbf{e}^T \mathbf{A} \mathbf{e}$ .

Using  $p_{i,j:\lambda}$ , we can write

$$\lambda(\lambda-1)u_3(\Phi(x), \Phi(y))p(x)p(y) = \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \max(p_{k,l:\lambda}(x, y), p_{l,k:\lambda}(x, y)) .$$

The equality (22) is obtained by substituting the equality

$$\begin{aligned} & \lambda(\lambda-1)\mathbb{E}[u_3(\Phi(\mathbf{Z}_e), \Phi(\tilde{\mathbf{Z}}_e))Z_e \tilde{Z}_e] \\ &= \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \iint Z_e \tilde{Z}_e \max(p_{k,l:\lambda}(x, y), p_{l,k:\lambda}(x, y)) dx dy \\ &= \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \left( \iint_{x < y} x y p_{k,l:\lambda}(x, y) dx dy + \iint_{x \geq y} x y p_{l,k:\lambda}(x, y) dx dy \right) \\ &= \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \left( \iint x y p_{k,l:\lambda}(x, y) dx dy + \iint x y p_{l,k:\lambda}(x, y) dx dy \right) \\ &= 2 \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \iint x y p_{k,l:\lambda}(x, y) dx dy \\ &= 2 \sum_{k=1}^{\lambda-1} \sum_{l=k+1}^{\lambda} w_k w_l \mathbb{E}[\mathcal{N}_{k:\lambda} \mathcal{N}_{l:\lambda}] \end{aligned}$$

into  $\mathbb{E}[u_3(\Phi(\mathbf{Z}_e), \Phi(\tilde{\mathbf{Z}}_e))Z^T \mathbf{A} \tilde{\mathbf{Z}}] = \mathbb{E}[u_3(\Phi(\mathbf{Z}_e), \Phi(\tilde{\mathbf{Z}}_e))Z_e \tilde{Z}_e] \mathbf{e}^T \mathbf{A} \mathbf{e}$ . The last equality is obtained by using the expression  $Z^T \mathbf{A} \tilde{\mathbf{Z}} = Z_e \tilde{Z}_e \mathbf{e}^T \mathbf{A} \mathbf{e} + Z_e \mathbf{e}^T \mathbf{A} \tilde{\mathbf{Z}}_{\perp} + \tilde{Z}_e \mathbf{e}^T \mathbf{A} \mathbf{Z}_{\perp} + Z_{\perp}^T \mathbf{A} \tilde{\mathbf{Z}}_{\perp}$ , the mutual independence between  $Z_e$ ,  $\tilde{Z}_e$ ,  $Z_{\perp}$ , and  $\tilde{\mathbf{Z}}_{\perp}$ , and the equalities  $\mathbb{E}[Z_{\perp}] = \mathbb{E}[\tilde{\mathbf{Z}}_{\perp}] = 0$ .  $\square$

### References

- [1] Y. Akimoto, A. Auger, N. Hansen, Quality gain analysis of the weighted recombination evolution strategy on general convex quadratic functions, in: Foundations of Genetic Algorithms, FOGA XIV, 2017, pp. 111–126.
- [2] N. Hansen, S. Kern, Evaluating the cma evolution strategy on multimodal test functions, in: Parallel Problem Solving from Nature, PPSN VIII, 2004, pp. 282–291.
- [3] N. Hansen, A. Auger, Principled design of continuous stochastic search: from theory to practice, in: Y. Borenstein, A. Moraglio (Eds.), Theory and Principled Methods for the Design of Metaheuristics, Springer, 2014.
- [4] N. Hansen, Invariance, self-adaptation and correlated mutations in evolution strategies, in: M. Schoenauer, K. Deb, G. Rudolph, X. Yao, E. Lutton, J.J.M. Guervós, H.-P. Schwefel (Eds.), Parallel Problem Solving from Nature, PPSN VI, Springer, 2000, pp. 355–364.
- [5] N. Hansen, A. Auger, R. Ros, S. Finck, P. Pošík, Comparing results of 31 algorithms from the black-box optimization benchmarking bbob-2009, in: Proceedings of Genetic and Evolutionary Computation Conference, 2010, pp. 1689–1696.
- [6] L.M. Rios, N.V. Sahinidis, Derivative-free optimization: a review of algorithms and comparison of software implementations, J. Global Optim. 56 (3) (2013) 1247–1293.
- [7] N. Hansen, A. Atamna, A. Auger, How to assess step-size adaptation mechanisms in randomised search, in: Parallel Problem Solving from Nature, PPSN XIII, Springer, 2014, pp. 60–69.
- [8] O. Krause, T. Glasmachers, C. Igel, Qualitative and quantitative assessment of step size adaptation rules, in: Proceedings of the 14th ACM/SIGEVO Conference on Foundations of Genetic Algorithms, FOGA '17, ACM, New York, NY, USA, 2017, pp. 139–148.
- [9] D.V. Arnold, Optimal weighted recombination, in: Foundations of Genetic Algorithms, FOGA VIII, Springer, 2005, pp. 215–237.
- [10] N. Hansen, A.S.P. Niederberger, L. Guzzella, P. Koumoutsakos, A method for handling uncertainty in evolutionary optimization with an application to feedback control of combustion, IEEE Trans. Evol. Comput. 13 (1) (2009) 180–197.

- [11] T. Yamaguchi, Y. Akimoto, Benchmarking the novel CMA-ES restart strategy using the search history on the bbob noiseless testbed, in: Proceedings of GECCO '17 Companion, 2017, pp. 1780–1787.
- [12] Y. Akimoto, N. Hansen, Online model selection for restricted covariance matrix adaptation, in: Parallel Problem Solving from Nature, PPSN XIV, 2016, pp. 3–13.
- [13] I. Rechenberg, *Evolutionstrategie '94*, Frommann-Holzboog, Stuttgart-Bad Cannstatt, 1994.
- [14] H.-G. Beyer, Towards a theory of 'evolution strategies': results for  $(1+, \lambda)$ -strategies on (nearly) arbitrary fitness functions, in: Parallel Problem Solving from Nature, PPSN III, 1994, pp. 58–67.
- [15] H.-G. Beyer, *The Theory of Evolution Strategies*, Nat. Comput. Ser., Springer-Verlag, 2001.
- [16] A. Auger, Convergence results for the  $(1, \lambda)$ -SA-ES using the theory of  $\varphi$ -irreducible markov chains, *Theoret. Comput. Sci.* 334 (1–3) (2005) 35–69.
- [17] A. Auger, N. Hansen, Reconsidering the progress rate theory for evolution strategies in finite dimensions, in: Proceedings of Genetic and Evolutionary Computation Conference, GECCO '06, 2006, pp. 445–452.
- [18] M. Jebalia, A. Auger, P. Liardet, Log-linear convergence and optimal bounds for the  $(1+1)$ -es, in: Evolution Artificielle, EA '07, 2008, pp. 207–218.
- [19] M. Jebalia, A. Auger, Log-linear convergence of the scale-invariant  $(\mu/\mu_w, \lambda)$ -ES and optimal  $\mu$  for intermediate recombination for large population sizes, in: Parallel Problem Solving from Nature, PPSN XI, 2010, pp. 52–62.
- [20] A. Auger, *Analysis of comparison-based stochastic continuous black-box optimization algorithms*, Habilitation, Université Paris-Sud, 2015.
- [21] R. Ros, N. Hansen, A simple modification in CMA-ES achieving linear time and space complexity, in: Parallel Problem Solving from Nature, PPSN X, 2008, pp. 296–305.
- [22] I. Loshchilov, A computationally efficient limited memory CMA-ES for large scale optimization, in: Proceedings of Genetic and Evolutionary Computation Conference, GECCO '14, 2014, pp. 397–404.
- [23] Y. Akimoto, N. Hansen, Projection-based restricted covariance matrix adaptation for high dimension, in: Proceedings of Genetic and Evolutionary Computation Conference, GECCO '16, 2016, pp. 197–204.
- [24] D.V. Arnold, On the use of evolution strategies for optimising certain positive definite quadratic forms, in: Proceedings of Genetic and Evolutionary Computation Conference, GECCO '07, 2007, pp. 634–641.
- [25] J. Jägersküpfer, How the  $(1+1)$  ES using isotropic mutations minimizes positive definite quadratic forms, *Theoret. Comput. Sci.* 361 (1) (2006) 38–56.
- [26] S. Finck, H.-G. Beyer, Weighted recombination evolution strategy on a class of pdqfs, in: Foundations of Genetic Algorithms, FOGA X, 2009, pp. 1–12.
- [27] H.-G. Beyer, A. Melkozerov, The dynamics of self-adaptive multirecombinant evolution strategies on the general ellipsoid model, *IEEE Trans. Evol. Comput.* 18 (5) (2014) 764–778.
- [28] H.-G. Beyer, M. Hellwig, The dynamics of cumulative step size adaptation on the ellipsoid model, *Evol. Comput.* 24 (1) (2016) 25–57.
- [29] A. Auger, D. Brockhoff, N. Hansen, Mirrored sampling in evolution strategies with weighted recombination, in: Proceedings of Genetic and Evolutionary Computation Conference, GECCO '11, 2011, pp. 861–868.
- [30] M. Jebalia, A. Auger, Log-Linear Convergence of the Scale-Invariant  $(\mu/\mu_w, \lambda)$ -ES and Optimal  $\mu$  for Intermediate Recombination for Large Population, Research Report RR-7275, INRIA, 2010.
- [31] O. Teytaud, S. Gelly, General lower bounds for evolutionary algorithms, in: Parallel Problem Solving from Nature, PPSN IX, 2006, pp. 21–31.
- [32] B. Laurent, P. Massart, Adaptive estimation of a quadratic functional by model selection, *Ann. Statist.* (2000) 1302–1338.
- [33] L. Devroye, *Non-Uniform Random Variate Generation*, Springer, New York, 1986.
- [34] H.-G. Beyer, Mutate large, but inherit small! On the analysis of rescaled mutations in  $(\tilde{1}, \tilde{\lambda})$ -ES with noisy fitness data, in: Parallel Problem Solving from Nature, PPSN V, 1998, pp. 109–118.
- [35] D.V. Arnold, Weighted multirecombination evolution strategies, *Theoret. Comput. Sci.* 361 (2006) 18–37.
- [36] N. Hansen, A. Ostermeier, Completely derandomized self-adaptation in evolution strategies, *Evol. Comput.* 9 (2) (2001) 159–195.
- [37] A. DasGupta, *Asymptotic Theory of Statistics and Probability*, Springer Science & Business Media, 2008.
- [38] H. Robbins, A remark on stirling's formula, *Amer. Math. Monthly* 62 (1) (1955) 26–29.