

# A New Look at the Covariance Matrix Estimation in Evolution Strategies

Silja Meyer-Nieberg<sup>(✉)</sup> and E. Kropat

Universität der Bundeswehr München,  
Werner-Heisenberg-Weg 39, 85577 Neubiberg, Germany  
{silja.meyer-nieberg,erik.kropat}@unibw.de

**Abstract.** Evolution strategies belong to the best performing modern natural computing methods for continuous optimization. This paper takes a new look at the covariance matrix adaptation, a mechanism which is central to the algorithm. The adaptation focusses strongly on the sample covariance. However, as known from modern statistics, this estimate may be of poor quality if certain conditions are not fulfilled. Unfortunately, this is often the case in practice. This paper compares the established methods for the covariance correction in evolution strategies with the approaches in modern statistics. Furthermore, it introduces and evaluates new covariance correction schemes.

**Keywords:** Evolutionary algorithms · Continuous optimization · Evolution strategies · Covariance · Shrinkage

## 1 Introduction

Black-Box optimization is an important subcategory of optimization. Over the years, several methods have been developed - ranging from simple pattern search over mesh adaptive methods to natural computing, see e.g. [1, 8, 10]. This paper focuses on evolution strategies (ESs) which represent well-performing meta-heuristics for continuous, non-linear optimization. In recent workshops on black-box optimization, see e.g. [15], variants of this particular subtype of evolutionary algorithms have emerged as one the best performing methods among a broad range of competitors stemming from natural computing. Evolution strategies rely primarily on random changes to move through the search space. These random changes, usually normally distributed random variables, must be controlled by adapting both, the extend and the direction of the movements.

Modern evolution strategies apply therefore covariance matrix and step-size adaptation – with great success. However, most methods use the common estimate of the population covariance matrix as one component to guide the search. Here, there may be room for further improvement, especially with regard to common application cases of evolution strategies which usually concern optimization in high-dimensional search spaces. For efficiency reasons, the population size  $\lambda$ , that is, the number of candidate solutions, is kept below the search space dimensionality  $N$  and scales usually with  $\mathcal{O}(\log(N))$  or with  $\mathcal{O}(N)$ . In other words,

either  $\lambda \ll N$  or  $\lambda \approx N$  which may represent a problem when using the sample covariance matrix. This even more so, since the sample size used in the estimation is just a fraction of the population size. Furthermore, the result is not robust against outliers which may appear in practical optimization which has often to cope with noise. This paper introduces and explores new approaches addressing the first problem by developing a new estimate for the covariance matrix. To our knowledge, these estimators have not been applied to evolution strategies before.

The paper is structured as follows: First, evolution strategies are introduced and common ways to adapt the covariance matrix are described and explained. Afterwards, we point out a potential dangerous weakness of the traditionally used estimate of the population covariance. Candidates for better estimates are presented and described in the following section. We propose and investigate several approaches ranging from a transfer of shrinkage estimators over a maximum entropy covariance selection principle to a new combination of both approaches. The quality of the resulting algorithms is assessed in the experimental test section. Conclusions and possible further research directions constitute the last part of the paper.

## 1.1 Evolution Strategies

Evolutionary algorithms (EAs) [10] are population-based stochastic search and optimization algorithms including today genetic algorithms, genetic programming, (natural) evolution strategies, evolutionary programming, and differential evolution. As a rule, they require only weak preconditions on the function to be optimized. Therefore, they are applicable in cases when only point-wise function evaluations are possible.

An evolutionary algorithm starts with an initial population of candidate solutions. The individuals are either drawn randomly from the search space or are initialized according to previous information on good solutions. A subset of the parent population is chosen for the creation of the offspring. This process is termed parent selection. Creation normally consists of recombination and mutation. While recombination combines traits from two or more parents, mutation is an unary operator and is realized by random perturbations. After the offspring have been created, survivor selection is performed to determine the next parent population. Evolutionary algorithms differ in the representation of the solutions and in the realization of the selection, recombination, and mutation operators.

Evolution strategies (ESs) [20,22] are a variant of evolutionary algorithms that is predominantly applied in continuous search spaces. Evolution strategies are commonly notated as  $(\mu/\rho, \lambda)$ -ESs. The parameter  $\mu$  stands for the size of the parent population. In the case of recombination,  $\rho$  parents are chosen randomly and are combined for the recombination result. While other forms exist, recombination usually consists of determining the weighted mean of the parents [4]. The result is then mutated by adding a normally distributed random variable with zero mean and covariance matrix  $\sigma^2 \mathbf{C}$ . While there are ESs that operate without recombination, the mutation process is essential and can be seen as the main search operator. Afterwards, the individuals are evaluated using the

function to be optimized or a derived function which allows an easy ranking of the population. Only the rank of an individual is important for the selection.

There are two main types of evolution strategies: Evolution strategies with “plus”-selection and ESs with “comma”-selection. The first select the  $\mu$ -best offspring and parents as the next parent population, where ESs with “comma”-selection discard the old parent population completely and take only the best offspring. Methods for adapting the scale factor  $\sigma$  or the full covariance matrix have received a lot of attention (see [19]). The main approaches are described in the following section.

## 1.2 Covariance Matrix Adaptation

First, the update of the covariance matrix is addressed. In evolution strategies two types exist: one applied in the *covariance matrix adaptation evolution strategy* (CMA-ES) [14] which considers past information from the search and an alternative used by the *covariance matrix self-adaptation evolution strategy* (CMSA-ES) [5] which focusses more on the present population.

The covariance matrix update of the CMA-ES is explained first. The CMA-ES uses weighted intermediate recombination, in other words, it computes the weighted centroid of the  $\mu$  best individuals of the population. This mean  $\mathbf{m}^{(g)}$  is used for creating all offspring by adding a random vector drawn from a normal distribution with covariance matrix  $(\sigma^{(g)})^2 \mathbf{C}^{(g)}$ , i.e., the actual covariance matrix consists of a general scaling factor (or step-size or mutation strength) and the matrix denoting the directions. Following usual notation in evolution strategies this matrix  $\mathbf{C}^{(g)}$  will be referred to as *covariance matrix* in the following.

The basis for the CMA update is the common estimate of the covariance matrix using the newly created population. Instead of considering the whole population for deriving the estimates, though, it introduces a bias towards good search regions by taking only the  $\mu$  best individuals into account. Furthermore, it does not estimate the mean anew but uses the weighted mean  $\mathbf{m}^{(g)}$ . Following [14],

$$\mathbf{y}_{m:\lambda}^{(g+1)} := \frac{1}{\sigma^{(g)}} \left( \mathbf{x}_{m:\lambda}^{(g+1)} - \mathbf{m}^{(g)} \right) \quad (1)$$

are determined with  $\mathbf{x}_{m:\lambda}$  denoting the  $m$ th best of the  $\lambda$  particle according to the fitness ranking. The rank- $\mu$  update obtains the covariance matrix as

$$\mathbf{C}_{\mu}^{(g+1)} := \sum_{m=1}^{\mu} w_m \mathbf{y}_{m:\lambda}^{(g+1)} (\mathbf{y}_{m:\lambda}^{(g+1)})^T \quad (2)$$

To derive reliable estimates larger population sizes are usually necessary which is detrimental with regard to the algorithm’s speed. Therefore, past information, that is, past covariance matrices are usually also considered

$$\mathbf{C}^{(g+1)} := (1 - c_{\mu}) \mathbf{C}^{(g)} + c_{\mu} \mathbf{C}_{\mu}^{(g+1)} \quad (3)$$

with parameter  $0 \leq c_{\mu} \leq 1$  determining the effective time-horizon. In CMA-ESs, it has been found that an enhance of the general search direction in the

covariance matrix is usual beneficial. For this, the concepts of the *evolutionary path* and the *rank-one-update* are introduced. As its name already suggests, an evolutionary path considers the path in the search space the population has taken so far. The weighted means serve as representatives. Defining

$$\mathbf{v}^{(g+1)} := \frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}}$$

the evolutionary path reads

$$\mathbf{p}_c^{(g+1)} := (1 - c_c)\mathbf{p}_c^{(g)} + \sqrt{c_c(2 - c_c)\mu_{\text{eff}}}\left(\frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}}\right). \quad (4)$$

For details on the parameters, see e.g. [12]. The evolutionary path gives a general search direction that the ES has taken in the recent past. In order to bias the covariance matrix accordingly, the rank-one-update

$$\mathbf{C}_1^{(g+1)} := \mathbf{p}_c^{(g+1)}(\mathbf{p}_c^{(g+1)})^T \quad (5)$$

is performed and used as a further component of the covariance matrix. A normal distribution with covariance  $\mathbf{C}_1^{(g+1)}$  leads towards a one-dimensional distribution on the line defined by  $\mathbf{p}_c^{(g+1)}$ . With (5) and (3), the final covariance update of the CMA-ES reads

$$\mathbf{C}^{(g+1)} := (1 - c_1 - c_\mu)\mathbf{C}^{(g)} + c_1\mathbf{C}_1^{(g+1)} + c_\mu\mathbf{C}_\mu^{(g+1)}. \quad (6)$$

The CMA-ES is one of the most powerful evolution strategies. However, as pointed out in [5], its scaling behavior with the population size is not good. The alternative approach of the CMSA-ES [5] updates the covariance matrix differently. Considering again the definition (1), the covariance update is a convex combination of the old covariance and the population covariance, i.e., the rank- $\mu$  update

$$\mathbf{C}^{(g+1)} := \left(1 - \frac{1}{c_\tau}\right)\mathbf{C}^{(g)} + \frac{1}{c_\tau} \sum_{m=1}^{\mu} w_m \mathbf{y}_{m:\lambda}^{(g+1)} (\mathbf{y}_{m:\lambda}^{(g+1)})^T \quad (7)$$

with the weights usually set to  $w_m = 1/\mu$ . See [5] for information on the free parameter  $c_\tau$ .

### 1.3 Step-Size Adaptation

The CMA-ES uses the so-called *cumulative step-size adaptation* (CSA) to control the scaling parameter (also called *step-size*, *mutation strength* or *step-length*) [12]. To this end, the CSA determines again an evolutionary path by summing the movement of the population centers

$$\mathbf{p}_\sigma^{(g+1)} = (1 - c_\sigma)\mathbf{p}_\sigma^{(g)} + \sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}}\mathbf{C}^{(g)^{-\frac{1}{2}}}\frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} \quad (8)$$

eliminating the influence of the covariance matrix and the step length. For a detailed description of the parameters, see [12]. The length of the path in (8) is important. In the case of short path lengths, several movement of the centers counteract each other which is an indication that the step-size is too large and should be reduced. If on the other hand, the ES takes several consecutive steps in approximately the same direction, progress and algorithm speed would be improved, if larger changes were possible. Long path lengths, therefore, are an indicator for a required increase of the step length. Ideally, the CSA should result in uncorrelated steps.

After some calculations, see [12], the ideal situation is revealed as standard normally distributed steps, which leads to

$$\ln(\sigma^{(g+1)}) = \ln(\sigma^{(g)}) + \frac{c_\sigma}{d_\sigma} \left( \frac{\|\mathbf{P}_\sigma^{(g+1)}\| - \mu_{\chi_n}}{\mu_{\chi_n}} \right) \quad (9)$$

as the CSA-rule. The change is multiplicative in order to avoid numerical problems and results in non-negative scaling parameters. The parameter  $\mu_{\chi_n}$  in (9) stands for the mean of the  $\chi$ -distribution with  $n$  degrees of freedom. If a random variable follows a  $\chi_n^2$  distribution, its square root is  $\chi$ -distributed. The degrees of freedom coincide with the search space dimension. The CSA-rule works well in many application cases. It can be shown, however, that the original CSA encounter problems in large noise regimes resulting in a loss of step-size control and premature convergence. Therefore, uncertainty handling procedures and other safeguards are advisable.

An alternative approach for adapting the step-size is *self-adaptation* first introduced in [20] and developed further in [22]. It subjects the strategy parameters of the mutation to evolution. In other words, the scaling parameter or in its full form, the whole covariance matrix, undergoes recombination, mutation, and indirect selection processes. The working principle is based on an indirect stochastic linkage between good individuals and appropriate parameters: On average good parameters should lead to better offspring than too large or too small values or misleading directions. Although self-adaptation has been developed to adapt the whole covariance matrix, it is used nowadays mainly to adapt the step-size or a diagonal covariance matrix. In the case of the mutation strength, usually a log-normal distribution

$$\sigma_l^{(g)} = \sigma_{\text{base}} \exp(\tau \mathcal{N}(0, 1)) \quad (10)$$

is used for mutation. The parameter  $\tau$  is called the *learning rate* and is usually chosen to scale with  $1/\sqrt{2N}$ . The variable  $\sigma_{\text{base}}$  is either the parental scale factor or the result of recombination. For the step-size, it is possible to apply the same type of recombination as for the positions although different forms – for instance a multiplicative combination – could be used instead. The self-adaptation of the step-size is referred to as  *$\sigma$ -self-adaptation* ( $\sigma$ SA) in the remainder of this paper.

The newly created mutation strength is then directly used in the mutation of the offspring. If the resulting offspring is sufficiently good, the scale factor is

passed to the next generation. The baseline  $\sigma_{\text{base}}$  is either the mutation strength of the parent or if recombination is used the recombination result. Self-adaptation with recombination has been shown to be “robust” against noise [3] and is used in the CMSA-ES as update rule for the scaling factor. In [5] it was found that the CMSA-ES performs comparably to the CMA-ES for smaller populations but is less computational expensive for larger population sizes.

## 2 Concerning the Covariance Estimator

The covariance matrix  $\mathbf{C}_\mu$  which appears in (2) and (7) can be interpreted as the sample covariance matrix with sample size  $\mu$ . Two differences are present. The first using  $\mu$  instead of  $\mu-1$  can be explained by using the known mean instead of an estimate. The second lies in the non-identically distributed random variables of the population since order statistics appear. We will disregard that problem for the time being.

In the case of identically independently distributed random variables, the estimate converges almost surely towards the “true” covariance  $\Sigma$  for  $\mu \rightarrow \infty$ . In addition, the sample covariance matrix is related (in our case equal) to the maximum likelihood (ML) estimator of  $\Sigma$ . Both facts serve a justification to take  $\mathbf{C}_\mu$  as the substitute for the unknown true covariance for large  $\mu$ . However, the quality of the estimate can be quite poor if  $\mu < N$  or even  $\mu \approx N$ .

This was first discovered by Stein [23, 24]. Stein’s phenomenon states that while the ML estimate is often seen as the best possible guess, its quality may be poor and can be improved in many cases. This holds especially for high-dimensional spaces. The same problem transfers to covariance matrix estimation, see [21]. Also recognized by Stein, in case of small ratios  $\mu/N$  the eigenstructure of  $\mathbf{C}_\mu$  may not agree well with the true eigenstructure of  $\Sigma$ . As stated in [17], the largest eigenvalue has a tendency towards too large values, whereas the smallest shows the opposite behavior. This results in a larger spectrum of the sample covariance matrix with respect to the true covariance for  $N/\mu \not\rightarrow 0$  for  $\mu, N \rightarrow \infty$  [2]. As found by Huber [16], a heavy tail distribution leads also to a distortion of the sample covariance.

In statistics, considerable efforts have been made to find more reliable and robust estimates. Owing to the great importance of the covariance matrix in data mining and other statistical analyses, work is still ongoing. The following section provides a short introduction before focussing on the approach used for evolution strategies.

## 3 Approaches for Estimating the Covariance

As stated above, the estimation of high-dimensional covariance matrices has received a lot of attention, see e.g. [6]. Several types have been introduced, for example: shrinkage estimators, banding and tapering estimators, sparse matrix transform estimators, and the graphical Lasso estimator. This paper concentrates on shrinkage estimators and on an idea inspired by a maximum entropy

approach. Both classes can be computed comparatively efficiently. Future research will consider other classes of estimators.

### 3.1 Shrinkage Estimators

Most (linear) shrinkage estimators use the convex combination

$$\mathbf{S}_{\text{est}}(\rho) = \rho \mathbf{F} + (1 - \rho) \mathbf{C}_\mu \quad (11)$$

with  $\mathbf{F}$  the *target* to correct the estimate provided by the sample covariance. The parameter  $\rho \in ]0, 1[$  is called the *shrinkage intensity*. Equation (11) is used to shrink the eigenvalues of  $\mathbf{C}_\mu$  towards the eigenvalues of  $\mathbf{F}$ . The shrinkage intensity  $\rho$  should be chosen to minimize

$$\mathbb{E} \left( \|\mathbf{S}_{\text{est}}(\rho) - \Sigma\|_F^2 \right) \quad (12)$$

with  $\|\cdot\|_F^2$  denoting the squared Frobenius norm with

$$\|\mathbf{A}\|_F^2 = \frac{1}{N} \text{Tr}[\mathbf{A}\mathbf{A}^T], \quad (13)$$

see [17]. To solve this problem, knowledge of the true covariance  $\Sigma$  would be required which is usually unobtainable.

Starting from (12), Ledoit and Wolf obtained an analytical expression for the optimal shrinkage intensity for the target  $\mathbf{F} = \text{Tr}(\mathbf{C}_\mu)/N \mathbf{I}$ . The result does not make assumptions on the underlying distribution. In the case of  $\mu \approx N$  or vastly different eigenvalues, the shrinkage estimator does not differ much from the sample covariance matrix, however.

Other authors introduced different estimators, see e.g. [7] or [6]). Ledoit and Wolfe themselves considered non-linear shrinkage estimators [18]. Most of the approaches require larger computational efforts. In the case of the non-linear shrinkage, for example, the authors are faced with a non-linear, non-convex optimization problem, which they solve by using sequential linear programming [18]. A general analytical expression is unobtainable, however.

Shrinkage estimators and other estimators aside from the standard case have not been used in evolution strategies before. A literature review resulted in one application in the case of Gaussian based estimation of distribution algorithms albeit with quite a different goal [9]. There, the learning of the covariance matrix during the run lead to non positive definite matrices. A shrinkage procedure was applied to “repair” the covariance matrix towards the required structure. The authors used a similar approach as in [17] but made the shrinkage intensity adaptable.

Interestingly, (3), (6), and (7) of the ES algorithm can be interpreted as a special case of shrinkage. In the case of the CMSA-ES, for example, the estimate is shrunk towards the old covariance matrix. The shrinkage intensity is determined by

$$c_\tau = 1 + \frac{N(N+1)}{2\mu} \quad (14)$$

as  $\rho = 1 - 1/c_\tau$ . As long as the increase of  $\mu$  with the dimensionality  $N$  is below  $\mathcal{O}(N^2)$ , the coefficient (14) approaches infinity for  $N \rightarrow \infty$ . Since the contribution of the sample covariance to the new covariance in (7) is weighted with  $1/c_\tau$ , its influence fades out for increasing dimensions. It is the aim of the paper to investigate whether a further shrinkage can improve the result.

Transferring shrinkage estimators to ESs must take the situation in which the estimation occurs into account since it differs from the assumptions in statistical literature. The covariance matrix  $\Sigma = \mathbf{C}^{g-1}$  that was used to create the offspring is known. The sample is based on rank-based selection, however, which differs from the iid case usually considered. Only if there were no selection pressure, the sample  $\mathbf{x}_1, \dots, \mathbf{x}_\mu$  would represent normally distributed random variables. In this context, it is interesting to note that the argumentation in [12] with respect to the setting of the CMA-ES parameter argues to choose the parameter so that the distribution of the random variables remains unchained as long as no selection pressure occurs. In other words, if  $\mathbf{p}^{(g)} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)})$  then also  $\mathbf{p}^{(g+1)}$  for both evolution paths, (4) and (8). However, due to sampling and using the covariance estimate, larger deviations may occur. Applying shrinkage could improve the situation. However, the choice of the target remains. Most shrinkage approaches consider diagonal matrices as shrinkage targets. If we were following that approach, we could choose the matrix  $\mathbf{F} = \text{diag}(\mathbf{C}_\mu)$ . This would leave the diagonal elements of the sample covariance matrix unchanged decreasing only the off-diagonal entries. However, a shrinkage towards a diagonal does not appear to be a good idea for optimizing functions that are not oriented towards the coordinate system.

### 3.2 A Maximum Entropy Covariance Estimation

Therefore, we make use of another concept following [25]. Confronted with the problem of determining a reliable covariance matrix by combining a sample covariance matrix with a pooled variance matrix, the authors introduced a *maximum entropy covariance selection principle*. Since a combination of covariance matrices also appears in evolution strategies, a closer look at their approach is interesting. Defining a population matrix  $\mathbf{C}_p$  and the sample covariance matrix  $\mathbf{S}_i$ , the mixture

$$\mathbf{S}_{mix}(\eta) = \eta \mathbf{C}_p + (1 - \eta) \mathbf{S}_i \quad (15)$$

was considered. In departure from usual approaches, focus lay on the combination of the two matrixes that maximizes the entropy. To this end, the coordinate system was changed to the eigenspace of  $\mathbf{S}_{mix} = \mathbf{C}_p + (1 - \eta) \mathbf{S}_i$ . Let  $\mathbf{M}_S$  denote the (normalized) eigenvectors of the mixture matrix. The representations of  $\mathbf{C}_p$  and  $\mathbf{S}_i$  in this coordination system read

$$\begin{aligned} \Phi^C &= \mathbf{M}_S^T \mathbf{C}_p \mathbf{M}_S \\ \Phi^S &= \mathbf{M}_S^T \mathbf{S}_i \mathbf{M}_S. \end{aligned} \quad (16)$$



Both matrices are usually not diagonal. To construct the new estimate for the covariance matrix,

$$\begin{aligned} \Lambda^C &= \text{diag}(\Phi^C) \\ \Lambda^S &= \text{diag}(\Phi^S) \end{aligned} \tag{17}$$

were determined. By taking  $\lambda_i = \max(\lambda_i^C, \lambda_i^S)$ , a covariance matrix estimate could finally be constructed via  $\mathbf{M}_S \Lambda \mathbf{M}_S^T$ . The approach maximizes the possible contributions to the principal direction of the mixture matrix and is based on a maximum entropy derivation for the estimation.

### 3.3 New Covariance Estimators

This paper proposes a combination of a shrinkage estimator and the basis transformation introduced [25] for a use in evolution strategies. This paper focuses on the CMSA-ES. The aim is to switch towards a suitable coordinate system and then either to discard the contributions of the sample covariance that are not properly aligned or to shrink the off-diagonal components. Two choices for the mixture matrix represent themselves. The first

$$\mathbf{S}_{mix} = \mathbf{C}^g + \mathbf{C}_\mu \tag{18}$$

is chosen in accordance to [25]. The second takes the covariance result that would have been used in the original CMSA-ES

$$\mathbf{S}_{mix} = (1 - c_\tau) \mathbf{C}^g + c_\tau \mathbf{C}_\mu \tag{19}$$

and introduces a single step recursion which may be more appropriate for small population sizes. Both choices will be investigated in this paper. They in turn can be coupled with several further ways to proceed and to construct the new covariance matrix. Switching towards the eigenspace of  $\mathbf{S}_{mix}$ , results in the covariance matrix representations  $\Phi_\mu := \mathbf{M}_S^T \mathbf{C}_\mu \mathbf{M}_S$  and  $\Phi_\Sigma := \mathbf{M}_S^T \mathbf{C}^g \mathbf{M}_S$ .

1. The first approach for constructing a new estimate of the sample covariance is to apply the principle of maximal contribution to the axes from [25] and to determine

$$\Lambda_\mu = \max \left( \text{diag}(\Phi_\mu), \text{diag}(\Phi_\Sigma) \right) \tag{20}$$

The sample covariance matrix can then be computed as  $\mathbf{C}'_\mu = \mathbf{M}_S \Lambda_\mu \mathbf{M}_S^T$ .

2. Another approach would be to discard all entries of  $\Phi_\mu$  except the diagonal

$$\Lambda_\mu = \text{diag}(\Phi_\mu) \tag{21}$$

3. A third approach consists of applying a shrinkage estimator like

$$\Phi_\mu^S = (1 - \rho) \Phi_\mu + \rho \text{diag}(\Phi_\mu). \tag{22}$$

This approach does not discard the off-diagonal entries completely. The shrinkage intensity  $\rho$  remains to be determined.

## 4 Experimental Evaluation

This section describes the experiments that were performed to explore the new approaches. For our investigation, the CMSA-ES version is considered since it operates just with the population covariance matrix and effects from changing the estimate should be easier to discerned. The competitors consist of algorithms which use shrinkage estimators as defined in (18) to (22). This code is not optimized for performance with respect to absolute computing time, since this paper aims at a proof of concept. The experiments are performed for the search space dimensions  $N = 2, 5, 10,$  and  $20$ . The maximal number of fitness evaluations is  $FE_{\max} = 2 \times 10^4 N$ . The CMSA-ES versions use  $\lambda = \lfloor \log(3N) + 8 \rfloor$  offspring and  $\mu = \lceil \lambda/4 \rceil$  parents. The start position of the algorithms is randomly chosen from a normal distribution with mean zero and standard deviation of 0.5. A run terminates prematurely if the difference between the best value obtained so far and the optimal fitness value  $|f_{\text{best}} - f_{\text{opt}}|$  is below a predefined precision set to  $10^{-8}$ . For each fitness function and dimension, 15 runs are used.

### 4.1 Test Suite

The experiments are performed with the black box optimization benchmarking (BBOB) software framework and the test suite introduced for the black box optimization workshops, see [13]. The aim of the workshop is to benchmark and compare metaheuristics and other direct search methods for continuous optimization. The framework allows the plug-in of algorithms adhering to a common interface and provides a comfortable way of generating the results in form of tables and figures.

The test suite contains noisy and noise-less functions with the position of the optimum changing randomly from run to run. This paper focuses on the 24 noise-less functions [11]. They can be divided into four classes: separable functions (function ids 1–5), functions with low/moderate conditioning (ids 6–9), functions with high conditioning (ids 10–14), and two groups of multimodal functions (ids 15–24).

### 4.2 Performance Measure

The following performance measure is used in accordance to [13]. The expected running time (ERT) gives the expected value of the function evaluations ( $f$ -evaluations) the algorithm needs to reach the target value with the required precision for the first time, see [13]. In this paper, we use

$$\text{ERT} = \frac{\#(FEs(f_{\text{best}} \geq f_{\text{target}}))}{\#succ} \quad (23)$$

as an estimate by summing up the fitness evaluations  $FEs(f_{\text{best}} \geq f_{\text{target}})$  of each run until the fitness of the best individual is smaller than the target value, divided by all successful runs (Fig. 1).

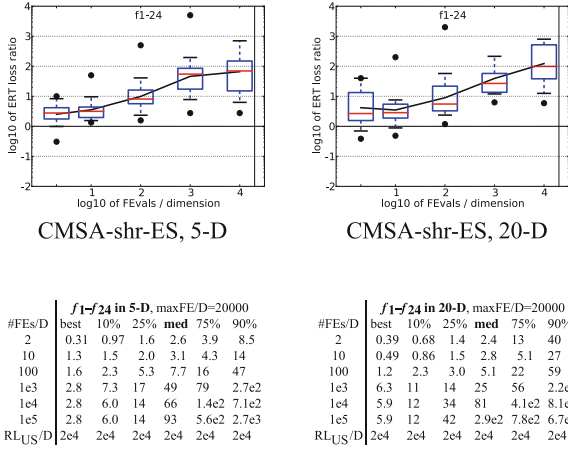
### 4.3 Results and Discussion

Due to space restrictions, Figure 3 and Table 1 and Fig. 2 show only the results from the best experiments which were achieved for the variant which used (22) together with (19) as the transformation matrix (called CMSA-shr-ES in the following). First of all, it should be noted that there is no significant advantage to either algorithm for the test suite functions. Table 1 and Fig. 2 show the ERT loss ratio with respect to the best result from the BBOB 2009 workshop for predefined budgets given in the first column. The median performance of both algorithms improves with the dimension until the budget of  $10^3$  – which is interesting. An increase of the budget goes along with a decreased performance which is less pronounced for the CMSA-shr-ES in the case of the larger dimensional space. This indicates that the CMSA-shr-ES may perform more favorable in larger search spaces as envisioned. Further experiments which a larger maximal number of fitness evaluations and larger dimensional spaces will be conducted which should shed more light on the behavior. Furthermore, the decrease in performance with the budget hints at a search stagnation probably due to convergence into local optima. Restart strategies may be beneficial, but since they have to be fitted to the algorithms, we do not apply them in the present paper.

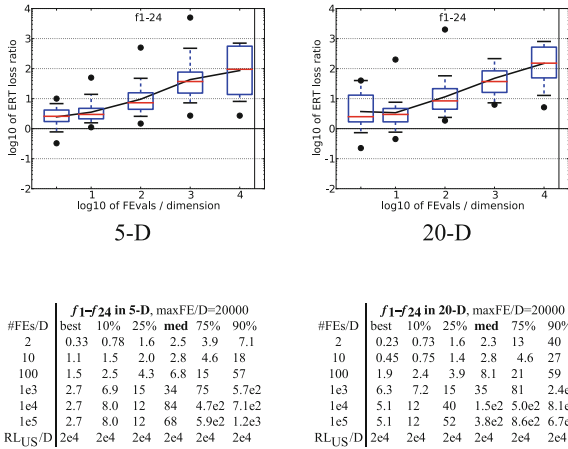
Figure 3 shows the expected running time for reaching the precision of  $10^{-8}$  for all 24 functions and search space dimensionalities. In the case of the separable functions (1–5), both algorithms show a very similar behavior, succeeding in optimizing the first two functions and exhibiting difficulties in the case of the difficult rastrigin variants. On the linear slope, the original CMSA-ES shows fewer expected function evaluations for smaller dimensions which starts to change when the dimensionality is increased. For the functions with ids 6–9, with moderate condition numbers, there are advantages to the CMSA-shr-ES, with the exception of the rotated rosenbrock (9). Most of the functions with high conditioning, ids 10–12, and 14, can be solved by both variants with slightly better results for the CMSA-ES. The sharp ridge (id 13) appears as problematic, with the CMSA-shr-ES showing fewer fitness evaluations for hitting the various precisions goals in Table 1.

Interestingly, the CMSA-shr variant seems to perform better for the difficult multimodal functions, e.g., Gallaghers 101 peak function, a finding which should be explored in more detail. The results for the last two multimodal functions can be explained in part in that the computing resources were insufficient for the optimization. Even the best performing algorithms from the BBOB workshop needed more resources than we used in our experiments.

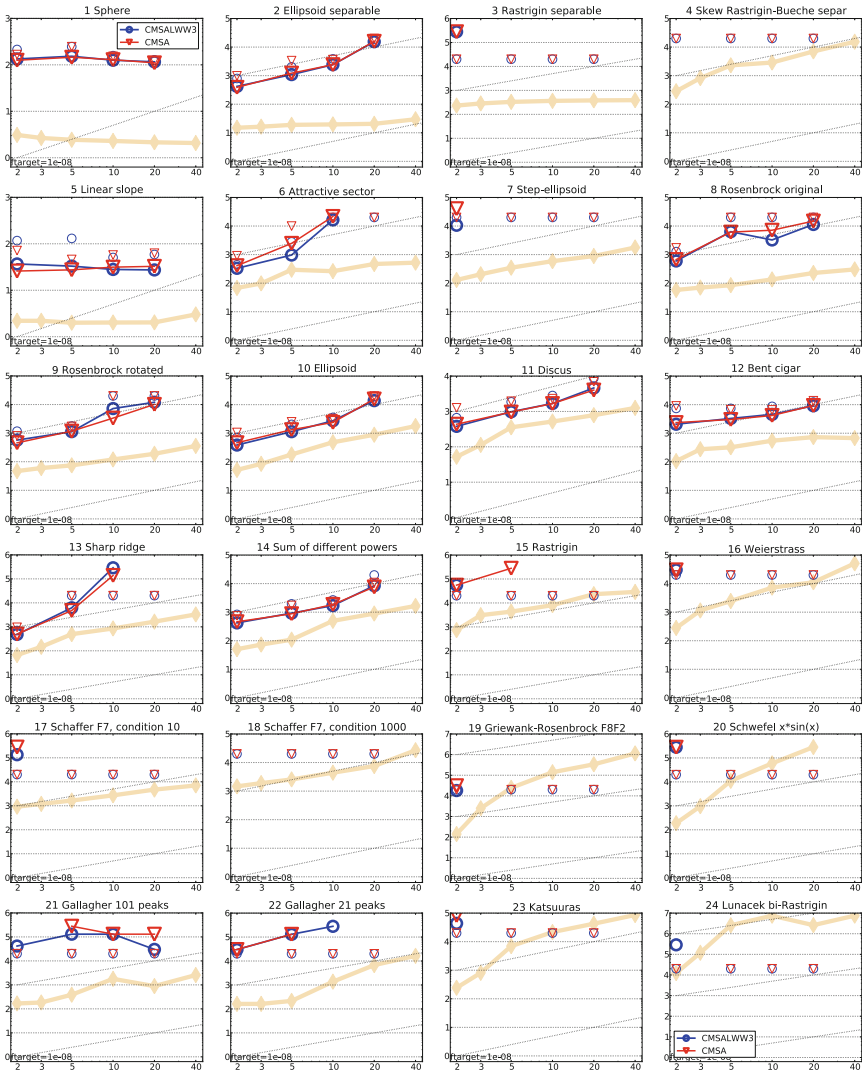
Further experiments will be conducted in order to shed more light on the behavior. Special attention will be given to the choice of the shrinkage factor, since its setting is unlikely to be optimal and may have influenced the outcome strongly. Furthermore, the question remains whether the population size should be increased for the self-adaptation process. Also, larger search space dimensionalities than  $N = 20$  are of interest.



**Fig. 1.** The CMSA-shr-ES. ERT loss ratio (in number of  $f$ -evaluations divided by dimension) divided by the best ERT seen in GECCO-BBOB-2009 for the target  $f_{target}$ , or, if the best algorithm reached a better target within the budget, the budget divided by the best ERT. Line: geometric mean. Box-Whisker error bar: 25–75 %-ile with median (box), 10–90 %-ile (caps), and minimum and maximum ERT loss ratio (points). The vertical line gives the maximal number of function evaluations in a single trial in this function subset.



**Fig. 2.** The CMSA-ES. ERT loss ratio (in number of  $f$ -evaluations divided by dimension) divided by the best ERT seen in GECCO-BBOB-2009 for the target  $f_{target}$ , or, if the best algorithm reached a better target within the budget, the budget divided by the best ERT. Line: geometric mean. Box-Whisker error bar: 25–75 %-ile with median (box), 10–90 %-ile (caps), and minimum and maximum ERT loss ratio (points). The vertical line gives the maximal number of function evaluations in a single trial in this function subset.



**Fig. 3.** Expected running time ERT in number of  $f$ -evaluations) divided by dimension for target function value as  $\log_{10}$  values versus dimension. Different symbols correspond to different algorithms given in the legend of  $f_1$  and  $f_{24}$ . Light symbols give the maximum number of function evaluations from the longest trial divided by dimension. Horizontal lines give linear scaling, slanted dotted lines give quadratic scaling. Black stars indicate statistically better result compared to all other algorithms with  $p < 0.01$  and Bonferroni correction number of dimensions (six). Legend: .1: CMSA-S is CMSA-shr-ES and 2: CMSA is CMSA-ES.



## 5 Conclusions

Evolution strategies are well performing variants of evolutionary algorithms used in continuous optimization. They utilize normally distributed mutations as their main search procedure. Their performance depends on the control of the mutation process which is governed by adapting step-sizes and covariance matrices. One possible improvement concerns the covariance matrix adaptation which makes use of the sample covariance matrix. In statistical research, this estimate has been identified as not agreeing well with the true covariance for the case of large dimensional spaces and small sample sizes, or more correctly for sample sizes that do not increase sufficiently fast with the dimensionality.

While modern approaches for covariance matrix adaptation correct the estimate, the question arises whether the performance of these evolutionary algorithms may be further improved by applying other estimators for the covariance.

This paper took a closer look at covariance estimation in evolution strategies and provided a comparison with approaches in modern statistics. Furthermore, it introduced and discussed new adaptation schemes for use in optimization. In cases, where the fitness function requires highly different eigenvalues and a rotation other than the cartesian coordinate system. Therefore, a switch towards the eigenspace of the covariance matrix was proposed in this paper and investigated in experiments on the BBOB test suite. While work remains to be done, this paper provided an important first step on the way.

## References

1. Audet, C.: A survey on direct search methods for blackbox optimization and their applications. In: Pardalos, P.M., Rassias, T.M. (eds.) *Mathematics without Boundaries: Surveys in Interdisciplinary Research*. Springer, New York (2013). Also *Les Cahiers du GERAD G-2012-53* (2012)
2. Bai, Z.D., Silverstein, J.W.: No eigenvalues outside the support of the limiting spectral distribution of large-dimensional sample covariance matrices. *Ann. Probab.* **26**(1), 316–345 (1998)
3. Beyer, H.G., Meyer-Nieberg, S.: Self-adaptation of evolution strategies under noisy fitness evaluations. *Genet. Program. Evolvable Mach.* **7**(4), 295–328 (2006)
4. Beyer, H.G., Schwefel, H.P.: Evolution strategies: a comprehensive introduction. *Nat. Comput.* **1**(1), 3–52 (2002)
5. Beyer, H.-G., Sendhoff, B.: Covariance matrix adaptation revisited – the CMSA evolution strategy –. In: Rudolph, G., Jansen, T., Lucas, S., Poloni, C., Beume, N. (eds.) *PPSN 2008*. LNCS, vol. 5199, pp. 123–132. Springer, Heidelberg (2008)
6. Chen, X., Wang, Z., McKeown, M.: Shrinkage-to-tapering estimation of large covariance matrices. *IEEE Trans. Signal Process.* **60**(11), 5640–5656 (2012)
7. Chen, Y., Wiesel, A., Eldar, Y.C., Hero, A.O.: Shrinkage algorithms for MMSE covariance estimation. *IEEE Trans. Signal Process.* **58**(10), 5016–5029 (2010)
8. Sarkar, M., Theuwissen, A.: Introduction. In: Sarkar, M., Theuwissen, A. (eds.) *A Biologically Inspired CMOS Image Sensor*. SCI, vol. 461, pp. 1–14. Springer, Heidelberg (2013)
9. Dong, W., Yao, X.: Covariance matrix repairing in gaussian based EDAs. In: *IEEE Congress on, Evolutionary Computation, CEC 2007*, pp. 415–422 (2007)

10. Eiben, A.E., Smith, J.E.: Introduction to Evolutionary Computing. Natural Computing Series. Springer, Berlin (2003)
11. Finck, S., Hansen, N., Ros, R., Auger, A.: Real-parameter black-box optimization benchmarking 2010: presentation of the noiseless functions. Technical report, Institute National de Recherche en Informatique et Automatique (2010) 2009/22
12. Hansen, N.: The CMA evolution strategy: a comparing review. In: Lozano J., et al. (eds.) Towards a new evolutionary computation. Advances in estimation of distribution algorithms, pp. 75–102. Springer (2006)
13. Hansen, N., Auger, A., Finck, S., Ros, R.: Real-parameter black-box optimization benchmarking 2012: experimental setup. Technical report, INRIA (2012). <http://coco.gforge.inria.fr/bbob2012-downloads>
14. Hansen, N., Ostermeier, A.: Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation* **9**(2), 159–195 (2001)
15. Hansen, N., Auger, A., Ros, R., Finck, S., Pošík, P.: Comparing results of 31 algorithms from the black-box optimization benchmarking BBOB-2009. In: Proceedings of the 12th Annual Conference Companion on Genetic and Evolutionary Computation, GECCO 2010, pp. 1689–1696. ACM, New York (2010). <http://doi.acm.org/10.1145/1830761.1830790>
16. Huber, P.J.: Robust Statistics. Wiley, New York (1981)
17. Ledoit, O., Wolf, M.: A well-conditioned estimator for large dimensional covariance matrices. *J. Multivar. Anal. Arch.* **88**(2), 265–411 (2004)
18. Ledoit, O., Wolf, M.: Non-linear shrinkage estimation of large dimensional covariance matrices. *Ann. Stat.* **40**(2), 1024–1060 (2012)
19. Meyer-Nieberg, S., Beyer, H.G.: Self-adaptation in evolutionary algorithms. In: Lobo, F., Lima, C., Michalewicz, Z. (eds.) Parameter Setting in Evolutionary Algorithms, pp. 47–76. Springer Verlag, Heidelberg (2007)
20. Rechenberg, I.: Evolutionsstrategie: Optimierung technischer Systeme nach Prinzipien der biologischen Evolution. Frommann-Holzboog Verlag, Stuttgart (1973)
21. Schäffer, J., Strimmer, K.: A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Stat. Appl. Genet. Mol. Biol.* **4**(1), Article 32 (2005)
22. Schwefel, H.P.: Numerical Optimization of Computer Models. Wiley, Chichester (1981)
23. Stein, C.: Inadmissibility of the usual estimator for the mean of a multivariate distribution. In: Proceedings 3rd Berkeley Symposium Mathematical Statistics and Probability, vol. 1, pp. 197–206. Berkeley, CA (1956)
24. Stein, C.: Estimation of a covariance matrix. In: Rietz Lecture, 39th Annual Meeting. IMS, Atlanta, GA (1975)
25. Thomaz, C.E., Gillies, D., Feitosa, R.: A new covariance estimate for bayesian classifiers in biometric recognition. *IEEE Trans. Circuits Syst. Video Technol.* **14**(2), 214–223 (2004)