

The CMA Evolution Strategy: A Tutorial

Nikolaus Hansen

January 18, 2009

Contents

Nomenclature	2
0 Preliminaries	3
0.1 Eigendecomposition of a Positive Definite Matrix	4
0.2 The Multivariate Normal Distribution	5
0.3 Randomized Black Box Optimization	6
0.4 Hessian and Covariance Matrices	7
1 Basic Equation: Sampling	7
2 Selection and Recombination: Moving the Mean	8
3 Adapting the Covariance Matrix	9
3.1 Estimating the Covariance Matrix From Scratch	9
3.2 Rank- μ -Update	11
3.3 Rank-One-Update	13
3.3.1 A Different Viewpoint	13
3.3.2 Cumulation: Utilizing the Evolution Path	15
3.4 Combining Rank- μ -Update and Cumulation	16
4 Step-Size Control	16
5 Discussion	20
A Algorithm Summary: The CMA-ES	24
A.1 Implementational Concerns	26
B MATLAB Source Code	29
C Reformulation of Learning Parameter c_{cov}	31

Nomenclature

We adopt the usual vector notation, where bold letters, \mathbf{v} , are column vectors, capital bold letters, \mathbf{A} , are matrices, and a transpose is denoted by \mathbf{v}^T . A list of used abbreviations and symbols is given in alphabetical order.

Abbreviations

CMA Covariance Matrix Adaptation

EMNA Estimation of Multivariate Normal Algorithm

ES Evolution Strategy

$(\mu/\mu_{\{I,W\}}, \lambda)$ -ES, Evolution Strategy with μ parents, with recombination of all μ parents, either Intermediate or Weighted, and λ offspring.

RHS Right Hand Side.

Greek symbols

$\lambda \geq 2$, population size, sample size, number of offspring, see (5).

$\mu \leq \lambda$ parent number, number of selected search points in the population, see (6).

μ_{cov} , parameter for weighting between rank-one and rank- μ update, see (27).

$\mu_{\text{eff}} = (\sum_{i=1}^{\mu} w_i^2)^{-1} = 1/\|\mathbf{w}\|^2$, the variance effective selection mass, see (8).

$\sigma^{(g)} \in \mathbb{R}_+$, step-size.

Latin symbols

$\mathbf{B} \in \mathbb{R}^n$, an orthogonal matrix. Columns of \mathbf{B} are eigenvectors of \mathbf{C} with unit length and correspond to the diagonal elements of \mathbf{D} .

$\mathbf{C}^{(g)} \in \mathbb{R}^{n \times n}$, covariance matrix at generation g .

c_{ii} , diagonal elements of \mathbf{C} .

$c_c \leq 1$, learning rate for cumulation for the rank-one update of the covariance matrix, see (22) and (42), and Table 1.

$c_{\text{cov}} \leq 1$, learning rate for the covariance matrix update, see (14), (26), (27), and (43), and Table 1.

$c_\sigma < 1$, learning rate for the cumulation for the step-size control, see (28) and (40), and Table 1.

$\mathbf{D} \in \mathbb{R}^n$, a diagonal matrix. The diagonal elements of \mathbf{D} are square roots of eigenvalues of \mathbf{C} and correspond to the respective columns of \mathbf{B} .

$d_i > 0$, diagonal elements of diagonal matrix \mathbf{D} , d_i^2 are eigenvalues of \mathbf{C} .

$d_\sigma \approx 1$, damping parameter for step-size update, see (29), (34), and (41).

E Expectation value

$f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, objective function (fitness function) to be minimized.

$f_{\text{sphere}} : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto \|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \sum_{i=1}^n x_i^2$.

$g \in \mathbb{N}_0$, generation counter, iteration number.

$\mathbf{I} \in \mathbb{R}^{n \times n}$, Identity matrix, unity matrix.

$\mathbf{m}^{(g)} \in \mathbb{R}^n$, mean value of the search distribution at generation g .

$n \in \mathbb{N}$, search space dimension, see f .

$\mathcal{N}(\mathbf{0}, \mathbf{I})$, multivariate normal distribution with zero mean and unity covariance matrix. A vector distributed according to $\mathcal{N}(\mathbf{0}, \mathbf{I})$ has independent, $(0, 1)$ -normally distributed components.

$\mathcal{N}(\mathbf{m}, \mathbf{C}) \sim \mathbf{m} + \mathcal{N}(\mathbf{0}, \mathbf{C})$, multivariate normal distribution with mean $\mathbf{m} \in \mathbb{R}^n$ and covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$. The matrix \mathbf{C} is symmetric and positive definite.

$\mathbf{p} \in \mathbb{R}^n$, evolution path, a sequence of successive (normalized) steps, the strategy takes over a number of generations.

w_i , where $i = 1, \dots, \mu$, recombination weights, see (6).

$\mathbf{x}_k^{(g+1)} \in \mathbb{R}^n$, k -th offspring/individual from generation $g + 1$. We also refer to $\mathbf{x}^{(g+1)}$, as search point, or object parameters/variables, commonly used synonyms are candidate solution, or design variables.

$\mathbf{x}_{i:\lambda}^{(g+1)}$, i -th best individual out of $\mathbf{x}_1^{(g+1)}, \dots, \mathbf{x}_\lambda^{(g+1)}$, see (5). The index $i : \lambda$ denotes the index of the i -th ranked individual and $f(\mathbf{x}_{1:\lambda}^{(g+1)}) \leq f(\mathbf{x}_{2:\lambda}^{(g+1)}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda}^{(g+1)})$, where f is the objective function to be minimized.

$\mathbf{y}_k^{(g+1)} = (\mathbf{x}_k^{(g+1)} - \mathbf{m}^{(g)}) / \sigma^{(g)}$ corresponding to $\mathbf{x}_k = \mathbf{m} + \sigma \mathbf{y}_k$.

0 Preliminaries

This tutorial introduces the CMA Evolution Strategy (ES), where CMA stands for Covariance Matrix Adaptation. The CMA-ES is a stochastic method for real-parameter (continuous domain) optimization of non-linear, non-convex functions (see also Section 0.3 below).¹ We try to motivate and derive the algorithm from intuitive concepts and from requirements of non-linear, non-convex search in continuous domain. In order to refer to the described algorithm, also cite [8]. For finding a concise algorithm description go directly to Appendix A. The respective Matlab source code is given in Appendix B.

Before we start to introduce the algorithm in Sect. 1, a few required fundamentals are summed up.

¹While CMA variants for *multi-objective* optimization and *elitistic* variants have been proposed, this tutorial is solely dedicated to single objective optimization and to non-elitistic truncation selection, also referred to as comma-selection.

0.1 Eigendecomposition of a Positive Definite Matrix

A symmetric, positive definite matrix, $C \in \mathbb{R}^{n \times n}$, is characterized in that for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ holds $\mathbf{x}^T C \mathbf{x} > 0$. The matrix C has an orthonormal basis of eigenvectors, $B = [\mathbf{b}_1, \dots, \mathbf{b}_n]$, with corresponding eigenvalues, $d_1^2, \dots, d_n^2 > 0$.

That means for each \mathbf{b}_i holds

$$C \mathbf{b}_i = d_i^2 \mathbf{b}_i . \quad (1)$$

The important message from (1) is that *eigenvectors are not rotated* by C . This feature uniquely distinguishes eigenvectors. Because we assume the orthogonal eigenvectors to be of unit length, $\mathbf{b}_i^T \mathbf{b}_j = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$, and $B^T B = I$ (obviously this means $B^{-1} = B^T$, and it follows $B B^T = I$). An basis of eigenvectors is practical, because for any $\mathbf{v} \in \mathbb{R}^n$ we can find coefficients α_i , such that $\mathbf{v} = \sum_i \alpha_i \mathbf{b}_i$, and then we have $C \mathbf{v} = \sum_i d_i^2 \alpha_i \mathbf{b}_i$.

The eigendecomposition of C obeys

$$C = B D^2 B^T , \quad (2)$$

where

B is an orthogonal matrix, $B^T B = B B^T = I$. Columns of B form an orthonormal basis of eigenvectors.

$D^2 = B C B^T = \text{diag}(d_1^2, \dots, d_n^2)$ is a diagonal matrix with eigenvalues of C as diagonal elements.

$D = \text{diag}(d_1, \dots, d_n)$ is a diagonal matrix with square roots of eigenvalues of C as diagonal elements.

The matrix decomposition (2) is unique, apart from signs of columns of B and permutations of columns in B and D^2 respectively, given all eigenvalues are different.²

Given the eigendecomposition (2), the inverse C^{-1} can be computed via

$$\begin{aligned} C^{-1} &= (B D^2 B^T)^{-1} \\ &= B^{T^{-1}} D^{-2} B^{-1} \\ &= B D^{-2} B^T \\ &= B \text{diag} \left(\frac{1}{d_1^2}, \dots, \frac{1}{d_n^2} \right) B^T , \end{aligned}$$

From (2) we naturally define

$$C^{\frac{1}{2}} = B D B^T \quad (3)$$

and therefore

$$\begin{aligned} C^{-\frac{1}{2}} &= B D^{-1} B^T \\ &= B \text{diag} \left(\frac{1}{d_1}, \dots, \frac{1}{d_n} \right) B^T \end{aligned}$$

²Given m eigenvalues are equal, any orthonormal basis of their m -dimensional subspace can be used as column vectors. For $m > 1$ there are infinitely many such bases.

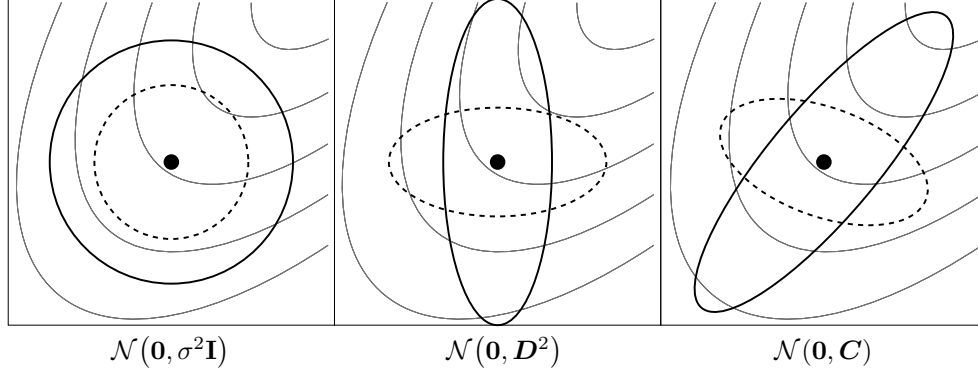


Figure 1: Ellipsoids depicting one- σ lines of equal density of six different normal distributions, where $\sigma \in \mathbb{R}_+$, \mathbf{D} is a diagonal matrix, and \mathbf{C} is a positive definite full covariance matrix. Thin lines depict possible objective function contour lines

0.2 The Multivariate Normal Distribution

A multivariate normal distribution, $\mathcal{N}(\mathbf{m}, \mathbf{C})$, has a unimodal, “bell-shaped” density, where the top of the bell (the modal value) corresponds to the distribution mean, \mathbf{m} . The distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is uniquely determined by its mean $\mathbf{m} \in \mathbb{R}^n$ and its symmetric and positive definite covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$. Covariance (positive definite) matrices have an appealing geometrical interpretation: they can be uniquely identified with the (hyper-)ellipsoid $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}^T \mathbf{C}^{-1} \mathbf{x} = 1\}$, as shown in Fig. 1. The ellipsoid is a surface of equal density of the distribution. The principal axes of the ellipsoid correspond to the eigenvectors of \mathbf{C} , the squared axes lengths correspond to the eigenvalues. The eigendecomposition is denoted by $\mathbf{C} = \mathbf{B}(\mathbf{D})^2 \mathbf{B}^T$ (see Sect. 0.1). If $\mathbf{D} = \sigma \mathbf{I}$, where $\sigma \in \mathbb{R}_+$ and \mathbf{I} denotes the identity matrix, $\mathbf{C} = \sigma^2 \mathbf{I}$ and the ellipsoid is isotropic (Fig. 1, left). If $\mathbf{B} = \mathbf{I}$, then $\mathbf{C} = \mathbf{D}^2$ is a diagonal matrix and the ellipsoid is axis parallel oriented (middle). In the coordinate system given by the columns of \mathbf{B} , the distribution $\mathcal{N}(\mathbf{0}, \mathbf{C})$ is always uncorrelated.

The normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ can be written in different ways.

$$\begin{aligned}
 \mathcal{N}(\mathbf{m}, \mathbf{C}) &\sim \mathbf{m} + \mathcal{N}(\mathbf{0}, \mathbf{C}) \\
 &\sim \mathbf{m} + \mathbf{C}^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
 &\sim \mathbf{m} + \mathbf{B} \underbrace{\mathbf{D} \mathbf{B}^T}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
 &\sim \mathbf{m} + \underbrace{\mathbf{B} \mathbf{D}}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{D}^2)} \mathcal{N}(\mathbf{0}, \mathbf{I}), \tag{4}
 \end{aligned}$$

where “ \sim ” denotes equality in distribution, and $\mathbf{C}^{\frac{1}{2}} = \mathbf{B} \mathbf{D} \mathbf{B}^T$. The last row can be well interpreted, from right to left

$\mathcal{N}(\mathbf{0}, \mathbf{I})$ produces an spherical (isotropic) distribution as in Fig. 1, left.

Initialize distribution parameters $\boldsymbol{\theta}^{(0)}$ For generation $g = 0, 1, 2, \dots$ Sample λ independent points from distribution $P(\mathbf{x} \boldsymbol{\theta}^{(g)}) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda$ Evaluate the sample $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$ on f Update parameters $\boldsymbol{\theta}^{(g+1)} = F_{\theta}(\boldsymbol{\theta}^{(g)}, (\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_\lambda, f(\mathbf{x}_\lambda)))$ break, if termination criterion met

Figure 2: Randomized black box search. $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the objective function

\mathbf{D} scales the spherical distribution within the coordinate axes as in Fig. 1, middle. $\mathcal{DN}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{D}^2)$ has n independent components. The matrix \mathbf{D} can be interpreted as (individual) step-size matrix and its diagonal entries are the standard deviations of the components.

\mathbf{B} defines a new orientation for the ellipsoid, where the new principal axes of the ellipsoid correspond to the columns of \mathbf{B} . Note that \mathbf{B} has $\frac{n^2-n}{2}$ degrees of freedom.

Equation (4) is useful to compute $\mathcal{N}(\mathbf{m}, \mathbf{C})$ distributed vectors, because $\mathcal{N}(\mathbf{0}, \mathbf{I})$ is a vector of independent $(0, 1)$ -normally distributed numbers that can easily be realized on a computer.

0.3 Randomized Black Box Optimization

We consider the black box search scenario, where we want to *minimize an objective function* (or *cost function* or *fitness function*)

$$\begin{aligned}
f & : \mathbb{R}^n \rightarrow \mathbb{R} \\
\mathbf{x} & \mapsto f(\mathbf{x}) .
\end{aligned}$$

The *objective* is to find one or more search points (candidate solutions), $\mathbf{x} \in \mathbb{R}^n$, with a function value, $f(\mathbf{x})$, as small as possible. We do not state the objective of searching for a global optimum, as this is often neither feasible nor relevant in practice. *Black box* optimization refers to the situation, where function values of evaluated search points are the only accessible information on f .³ The search points to be evaluated can be freely chosen. We define the *search costs* as the number of executed function evaluations, in other words the amount of information we needed to acquire from f .⁴ Any performance measure must consider the search costs *together* with the achieved objective function value.⁵

A randomized black box search algorithm is outlined in Fig. 2. In the CMA Evolution Strategy the search distribution, P , is a multivariate normal distribution. Given all variances and covariances, the normal distribution has the largest entropy of all distributions in \mathbb{R}^n .

³Knowledge about the underlying optimization problem might well enter the composition of f and the chosen problem *encoding*.

⁴Also f is sometimes denoted as *cost function*, but it should not to be confused with the *search costs*.

⁵A performance measure can be obtained from a number of trials as, for example, the mean number of function evaluations to reach a given function value, or the median best function value obtained after a given number of function evaluations.

Furthermore, coordinate directions are not distinguished in any way. Both makes the normal distribution a particularly attractive candidate for randomized search.

Randomized search algorithms are regarded to be robust in a rugged search landscape, which can comprise discontinuities, (sharp) ridges, or local optima. The covariance matrix adaptation (CMA) in particular is designed to tackle, additionally, ill-conditioned and non-separable⁶ problems.

0.4 Hessian and Covariance Matrices

We consider the convex-quadratic objective function $f_H : \mathbf{x} \mapsto \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x}$, where the Hessian matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ is a positive definite matrix. Given a search distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$, there is a close relation between \mathbf{H} and \mathbf{C} : Setting $\mathbf{C} = \mathbf{H}^{-1}$ on f_H is equivalent to optimizing the isotropic function $f_{\text{sphere}}(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{x} = \frac{1}{2} \sum_i x_i^2$ (where $\mathbf{H} = \mathbf{I}$) with $\mathbf{C} = \mathbf{I}$.⁷ That is, on convex-quadratic objective functions, setting the covariance matrix of the search distribution to the inverse Hessian matrix is equivalent to rescaling the ellipsoid function into a spherical one. Consequently, we assume that the optimal covariance matrix equals to the inverse Hessian matrix, up to a constant factor.⁸ Furthermore, choosing a covariance matrix or choosing a respective affine linear transformation of the search space (*i.e.* of \mathbf{x}) is equivalent [7], because for any full rank $n \times n$ -matrix \mathbf{A} we find a positive definite Hessian such that $\frac{1}{2} (\mathbf{A} \mathbf{x})^T \mathbf{A} \mathbf{x} = \frac{1}{2} \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x}$.

Consequently, the **objective** of covariance matrix adaptation is to *approximate the inverse Hessian matrix*, similar to a quasi-Newton method. More general, the objective is to suit the search distribution to the contour lines of the objective function f . In Fig. 1 the solid-line distribution in the right figure follows the objective function contours most suitably, and it is easy to foresee that it will aid to approach the optimum the most.

The **condition number** of a positive definite matrix \mathbf{A} is defined via the Euclidean norm: $\text{cond}(\mathbf{A}) \stackrel{\text{def}}{=} \|\mathbf{A}\| \times \|\mathbf{A}^{-1}\|$, where $\|\mathbf{A}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{A} \mathbf{x}\|$. For a positive definite (Hessian or covariance) matrix \mathbf{A} holds $\|\mathbf{A}\| = \lambda_{\max}$ and $\text{cond}(\mathbf{A}) = \frac{\lambda_{\max}}{\lambda_{\min}} \geq 1$, where λ_{\max} and λ_{\min} are the largest and smallest eigenvalue of \mathbf{A} .

1 Basic Equation: Sampling

In the CMA Evolution Strategy, a population of new search points (individuals, offspring) is generated by sampling a multivariate normal distribution.⁹ The basic equation for sampling the search points, for generation number $g = 0, 1, 2, \dots$, reads¹⁰

$$\mathbf{x}_k^{(g+1)} \sim \mathbf{m}^{(g)} + \sigma^{(g)} \mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)}) \quad \text{for } k = 1, \dots, \lambda \quad (5)$$

⁶An n -dimensional *separable* problem can be solved by solving n 1-dimensional problems separately, which is a far easier task.

⁷Also the initial mean value \mathbf{m} has to be transformed accordingly.

⁸Even though there is good intuition and strong empirical evidence for this statement, a rigorous proof is missing.

⁹Recall that, given all (co-)variances, the normal distribution has the largest entropy of all distributions in \mathbb{R}^n .

¹⁰Framed equations belong to the final algorithm of a CMA Evolution Strategy.

where

\sim denotes the same distribution on the left and right side.

$\mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)})$ is a multivariate normal distribution with zero mean and covariance matrix $\mathbf{C}^{(g)}$, see Sect. 0.2. It holds $\mathbf{m}^{(g)} + \sigma^{(g)}\mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)}) \sim \mathcal{N}(\mathbf{m}^{(g)}, (\sigma^{(g)})^2\mathbf{C}^{(g)})$.

$\mathbf{x}_k^{(g+1)} \in \mathbb{R}^n$, k -th offspring (individual, search point) from generation $g + 1$.

$\mathbf{m}^{(g)} \in \mathbb{R}^n$, mean value of the search distribution at generation g .

$\sigma^{(g)} \in \mathbb{R}_+$, “overall” standard deviation, step-size, at generation g .

$\mathbf{C}^{(g)} \in \mathbb{R}^{n \times n}$, covariance matrix at generation g . Up to the scalar factor $\sigma^{(g)2}$, $\mathbf{C}^{(g)}$ is the covariance matrix of the search distribution.

$\lambda \geq 2$, population size, sample size, number of offspring.

To define the complete iteration step, the remaining question is, how to calculate $\mathbf{m}^{(g+1)}$, $\mathbf{C}^{(g+1)}$, and $\sigma^{(g+1)}$ for the next generation $g + 1$. The next three sections will answer these questions, respectively. An algorithm summary with all parameter settings and MATLAB source code is given in Appendix A and B, respectively.

2 Selection and Recombination: Moving the Mean

The new mean $\mathbf{m}^{(g+1)}$ of the search distribution is a *weighted average of μ selected points* from the sample $\mathbf{x}_1^{(g+1)}, \dots, \mathbf{x}_\lambda^{(g+1)}$:

$$\mathbf{m}^{(g+1)} = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}^{(g+1)} \quad (6)$$

$$\sum_{i=1}^{\mu} w_i = 1, \quad w_1 \geq w_2 \geq \dots \geq w_\mu > 0 \quad (7)$$

where

$\mu \leq \lambda$ is the parent population size, *i.e.* the number of selected points.

$w_{i=1\dots\mu} \in \mathbb{R}_+$, positive weight coefficients for recombination. For $w_{i=1\dots\mu} = 1/\mu$, Equation (6) calculates the mean value of μ selected points.

$\mathbf{x}_{i:\lambda}^{(g+1)}$, i -th best individual out of $\mathbf{x}_1^{(g+1)}, \dots, \mathbf{x}_\lambda^{(g+1)}$ from (5). The index $i : \lambda$ denotes the index of the i -th ranked individual and $f(\mathbf{x}_{1:\lambda}^{(g+1)}) \leq f(\mathbf{x}_{2:\lambda}^{(g+1)}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda}^{(g+1)})$, where f is the objective function to be minimized.

Equation (6) implements *truncation selection* by choosing $\mu < \lambda$ out of λ offspring points. Assigning *different* weights w_i must also be interpreted as a selection mechanism. Equation (6) implements *weighted intermediate recombination* by taking $\mu > 1$ individuals into account for a weighted average.

The measure

$$\mu_{\text{eff}} = \left(\frac{\|\mathbf{w}\|_1}{\|\mathbf{w}\|_2} \right)^2 = \frac{1}{\|\mathbf{w}\|_2^2} = \left(\sum_{i=1}^{\mu} w_i^2 \right)^{-1} \quad (8)$$

will be repeatedly used in the following and can be paraphrased as *variance effective selection mass*. From the definition of w_i in (7) we derive $1 \leq \mu_{\text{eff}} \leq \mu$, and $\mu_{\text{eff}} = \mu$ for equal recombination weights, *i.e.* $w_i = 1/\mu$ for all $i = 1 \dots \mu$. Usually, $\mu_{\text{eff}} \approx \lambda/4$ indicates a reasonable setting of w_i . A typical setting could be $w_i \propto \mu - i + 1$, and $\mu \approx \lambda/2$.

3 Adapting the Covariance Matrix

In this section, the update of the covariance matrix, \mathbf{C} , is derived. We will start out estimating the covariance matrix from a single population of one generation (Sect. 3.1). For small populations this estimation is unreliable and an adaptation procedure has to be invented (rank- μ -update, Sect. 3.2). In the limit case only a single point can be used to update (adapt) the covariance matrix at each generation (rank-one-update, Sect. 3.3). The adaptation can be enhanced by exploiting dependencies between successive steps applying cumulation (Sect. 3.3.2). Finally we combine the rank- μ and rank-one updating methods (Sect. 3.4).

3.1 Estimating the Covariance Matrix From Scratch

For the moment we assume that the population contains enough information to reliably estimate a covariance matrix from the population.¹¹ For the sake of convenience we assume $\sigma^{(g)} = 1$ (see (5)) in this section. For $\sigma^{(g)} \neq 1$ the formulae hold except for a constant factor.

We can (re-)estimate the original covariance matrix $\mathbf{C}^{(g)}$ using the sampled population from (5), $\mathbf{x}_1^{(g+1)} \dots \mathbf{x}_\lambda^{(g+1)}$, via the empirical covariance matrix

$$\mathbf{C}_{\text{emp}}^{(g+1)} = \frac{1}{\lambda - 1} \sum_{i=1}^{\lambda} \left(\mathbf{x}_i^{(g+1)} - \frac{1}{\lambda} \sum_{j=1}^{\lambda} \mathbf{x}_j^{(g+1)} \right) \left(\mathbf{x}_i^{(g+1)} - \frac{1}{\lambda} \sum_{j=1}^{\lambda} \mathbf{x}_j^{(g+1)} \right)^{\text{T}}. \quad (9)$$

The empirical covariance matrix $\mathbf{C}_{\text{emp}}^{(g+1)}$ is an unbiased estimator of $\mathbf{C}^{(g)}$: assuming the $\mathbf{x}_i^{(g+1)}$, $i = 1 \dots \lambda$, to be random variables (rather than a realized sample), we have that $\mathbb{E}[\mathbf{C}_{\text{emp}}^{(g+1)} \mid \mathbf{C}^{(g)}] = \mathbf{C}^{(g)}$. Consider now a slightly different approach to get an estimator for $\mathbf{C}^{(g)}$.

$$\mathbf{C}_{\lambda}^{(g+1)} = \frac{1}{\lambda} \sum_{i=1}^{\lambda} \left(\mathbf{x}_i^{(g+1)} - \mathbf{m}^{(g)} \right) \left(\mathbf{x}_i^{(g+1)} - \mathbf{m}^{(g)} \right)^{\text{T}} \quad (10)$$

¹¹To re-estimate the covariance matrix, \mathbf{C} , from a $\mathcal{N}(\mathbf{0}, \mathbf{I})$ distributed sample such that $\text{cond}(\mathbf{C}) < 10$ a sample size $\lambda \geq 4n$ is needed, as can be observed in numerical experiments.

Also the matrix $\mathbf{C}_\lambda^{(g+1)}$ is an unbiased estimator of $\mathbf{C}^{(g)}$. The remarkable difference between (9) and (10) is the reference mean value. For $\mathbf{C}_{\text{emp}}^{(g+1)}$ it is the mean of the *actually realized* sample. For $\mathbf{C}_\lambda^{(g+1)}$ it is the *true* mean value, $\mathbf{m}^{(g)}$, of the sampled distribution (see (5)). Therefore, the estimators $\mathbf{C}_{\text{emp}}^{(g+1)}$ and $\mathbf{C}_\lambda^{(g+1)}$ can be interpreted differently: while $\mathbf{C}_{\text{emp}}^{(g+1)}$ estimates the distribution variance *within the sampled points*, $\mathbf{C}_\lambda^{(g+1)}$ estimates variances of sampled *steps*, $\mathbf{x}_i^{(g+1)} - \mathbf{m}^{(g)}$.

A minor difference between (9) and (10) is the different normalizations $\frac{1}{\lambda-1}$ versus $\frac{1}{\lambda}$, necessary to get an unbiased estimator in both cases. In (9) one degree of freedom is already taken by the inner summand. In order to get a maximum likelihood estimator in both cases $\frac{1}{\lambda}$ must be used.

Equation (10) re-estimates *the original* covariance matrix. To “estimate” a “better” covariance matrix (10) is modified and the same, *weighted selection* mechanism as in (6) is used [9].

$$\mathbf{C}_\mu^{(g+1)} = \sum_{i=1}^{\mu} w_i \left(\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)} \right) \left(\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)} \right)^{\text{T}} \quad (11)$$

The matrix $\mathbf{C}_\mu^{(g+1)}$ is an estimator for the distribution of *selected steps*, just as $\mathbf{C}_\lambda^{(g+1)}$ is an estimator of the original distribution of steps before selection. Sampling from $\mathbf{C}_\mu^{(g+1)}$ tends to reproduce selected, *i.e. successful* steps, giving a justification for what a “better” covariance matrix means.

Following [8], we compare (11) with the Estimation of Multivariate Normal Algorithm $\text{EMNA}_{\text{global}}$ [16, 17]. The covariance matrix in $\text{EMNA}_{\text{global}}$ reads, similar to (9),

$$\mathbf{C}_{\text{EMNA}_{\text{global}}}^{(g+1)} = \frac{1}{\mu} \sum_{i=1}^{\mu} \left(\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g+1)} \right) \left(\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g+1)} \right)^{\text{T}}, \quad (12)$$

where $\mathbf{m}^{(g+1)} = \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}_{i:\lambda}^{(g+1)}$. Similarly, applying the so-called Cross-Entropy method to continuous domain optimization [19] yields the covariance matrix $\frac{\mu}{\mu-1} \mathbf{C}_{\text{EMNA}_{\text{global}}}^{(g+1)}$, *i.e. the unbiased* empirical covariance matrix of the μ best points. In both cases the subtle, but most important difference to (11) is, again, the choice of the reference mean value.¹² Equation (12) estimates the variance *within* the selected population while (11) estimates selected steps. Equation (12) reveals always smaller variances than (11), because its reference mean value is the minimizer for the variances. Moreover, in most conceivable selection situations (12) decreases the variances compared to $\mathbf{C}^{(g)}$.

Figure 3 demonstrates the estimation results on a *linear* objective function for $\lambda = 150$, $\mu = 50$, and $w_i = 1/\mu$. Equation (11) geometrically increases the expected variance in direction of the gradient (where the selection takes place, here the diagonal), given ordinary settings for parent number μ and recombination weights w_1, \dots, w_μ . Equation (12) always decreases the variance in gradient direction geometrically fast! Therefore, (12) is highly susceptible to premature convergence, in particular with small parent populations, where the population cannot be expected to bracket the optimum at any time. However, for large values of μ in large populations with large initial variances, the impact of the different reference mean value can become marginal.

¹²Taking a weighted sum, $\sum_{i=1}^{\mu} w_i \dots$, instead of the mean, $\frac{1}{\mu} \sum_{i=1}^{\mu} \dots$, is an appealing, but less important, difference.

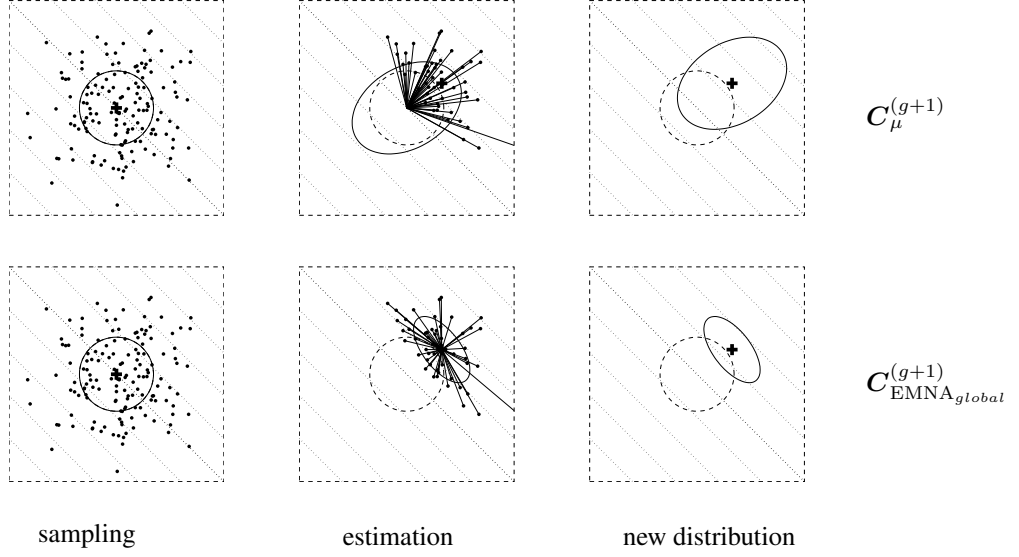


Figure 3: Estimation of the covariance matrix on $f_{\text{linear}}(\mathbf{x}) = -\sum_{i=1}^2 x_i$ to be minimized. Contour lines (*dotted*) indicate that the strategy should move toward the upper right corner. **Above:** estimation of $\mathbf{C}_{\mu}^{(g+1)}$ according to (11), where $w_i = 1/\mu$. **Below:** estimation of $\mathbf{C}_{\text{EMNA}_{\text{global}}}^{(g+1)}$ according to (12). Left: sample of $\lambda = 150 \mathcal{N}(\mathbf{0}, \mathbf{I})$ distributed points. Middle: the $\mu = 50$ selected points (*dots*) determining the entries for the estimation equation (*solid straight lines*). Right: search distribution of the next generation (*solid ellipsoids*). Given $w_i = 1/\mu$, estimation via $\mathbf{C}_{\mu}^{(g+1)}$ *increases* the expected variance in gradient direction for all $\mu < \lambda/2$, while estimation via $\mathbf{C}_{\text{EMNA}_{\text{global}}}^{(g+1)}$ *decreases* this variance for any $\mu < \lambda$ geometrically fast

In order to ensure with (5), (6), and (11), that $\mathbf{C}_{\mu}^{(g+1)}$ is a *reliable* estimator, the variance effective selection mass μ_{eff} (cf. (8)) must be large enough: getting condition numbers (cf. Sect. 0.4) smaller than ten for $\mathbf{C}_{\mu}^{(g)}$ on $f_{\text{sphere}}(\mathbf{x}) = \sum_{i=1}^n x_i^2$, requires $\mu_{\text{eff}} \approx 10n$. The next step is to circumvent this restriction on μ_{eff} .

3.2 Rank- μ -Update

To achieve *fast* search (opposite to *more robust* or *more global* search), e.g. competitive performance on f_{sphere} , the population size λ must be small. Because $\mu_{\text{eff}} \approx \lambda/4$ also μ_{eff} must be small and we may assume, e.g., $\mu_{\text{eff}} \leq 1 + \ln n$. Then, it is not possible to get a *reliable* estimator for a good covariance matrix from (11). As a remedy, information from previous generations is used additionally. For example, after a sufficient number of generations, the

mean of the estimated covariance matrices from all generations,

$$\mathbf{C}^{(g+1)} = \frac{1}{g+1} \sum_{i=0}^g \frac{1}{\sigma^{(i)2}} \mathbf{C}_\mu^{(i+1)} \quad (13)$$

becomes a reliable estimator for the selected steps. To make $\mathbf{C}_\mu^{(g)}$ from different generations comparable, the different $\sigma^{(i)}$ are incorporated. (Assuming $\sigma^{(i)} = 1$, (13) resembles the covariance matrix from the Estimation of Multivariate Normal Algorithm EMNA_i [17].)

In (13), all generation steps have the same weight. To assign recent generations a higher weight, exponential smoothing is introduced. Choosing $\mathbf{C}^{(0)} = \mathbf{I}$ to be the unity matrix and a learning rate $0 < c_{\text{cov}} \leq 1$, then $\mathbf{C}^{(g+1)}$ reads

$$\begin{aligned} \mathbf{C}^{(g+1)} &= (1 - c_{\text{cov}}) \mathbf{C}^{(g)} + c_{\text{cov}} \frac{1}{\sigma^{(g)2}} \mathbf{C}_\mu^{(g+1)} \\ &= (1 - c_{\text{cov}}) \mathbf{C}^{(g)} + c_{\text{cov}} \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}^{(g+1)} \mathbf{y}_{i:\lambda}^{(g+1)\text{T}} \end{aligned} \quad (14)$$

where

$c_{\text{cov}} \leq 1$ learning rate for updating the covariance matrix. For $c_{\text{cov}} = 1$, no prior information is retained and $\mathbf{C}^{(g+1)} = \frac{1}{\sigma^{(g)2}} \mathbf{C}_\mu^{(g+1)}$. For $c_{\text{cov}} = 0$, no learning takes place and $\mathbf{C}^{(g+1)} = \mathbf{C}^{(0)}$. Here, $c_{\text{cov}} \approx \min(1, \mu_{\text{eff}}/n^2)$ is a reasonable choice.

$$\mathbf{y}_{i:\lambda}^{(g+1)} = (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)}) / \sigma^{(g)}.$$

This covariance matrix update is called rank- μ -update [11], because the sum of outer products in (14) is of rank $\min(\mu, n)$ (with probability one). Note that this sum can even consist of a single term, if $\mu = 1$.

The number $1/c_{\text{cov}}$ is the *backward time horizon* that contains roughly 63% of the overall weight.

Because (14) expands to the weighted sum

$$\mathbf{C}^{(g+1)} = (1 - c_{\text{cov}})^{g+1} \mathbf{C}^{(0)} + c_{\text{cov}} \sum_{i=0}^g (1 - c_{\text{cov}})^{g-i} \frac{1}{\sigma^{(i)2}} \mathbf{C}_\mu^{(i+1)}, \quad (15)$$

the backward time horizon, Δg , where about 63% of the overall weight is summed up, is defined by

$$c_{\text{cov}} \sum_{i=g+1-\Delta g}^g (1 - c_{\text{cov}})^{g-i} \approx 0.63 \approx 1 - \frac{1}{e}. \quad (16)$$

Resolving the sum yields

$$(1 - c_{\text{cov}})^{\Delta g} \approx \frac{1}{e}, \quad (17)$$

and resolving for Δg , using the Taylor approximation for \ln , yields

$$\Delta g \approx \frac{1}{c_{\text{cov}}}. \quad (18)$$

That is, approximately 37% of the information in $\mathbf{C}^{(g+1)}$ is older than $1/c_{\text{cov}}$ generations, and, according to (17), the original weight is reduced by a factor of 0.37 after approximately $1/c_{\text{cov}}$ generations.¹³

The choice of c_{cov} is crucial. Small values lead to slow learning, too large values lead to a failure, because the covariance matrix degenerates. Fortunately, a good setting seems to be largely independent of the function to be optimized.¹⁴ A first order approximation for a good choice is $c_{\text{cov}} \approx \mu_{\text{eff}}/n^2$. Therefore, the characteristic time horizon for (14) is roughly n^2/μ_{eff} .

Even for the learning rate $c_{\text{cov}} = 1$, adapting the covariance matrix cannot be accomplished within one generation. The effect of the original sample distribution does not vanish until a sufficient number of generations. Assuming fixed search costs (number of function evaluations), a small population size λ allows a larger number of generations and therefore usually leads to a faster adaptation of the covariance matrix.

3.3 Rank-One-Update

In Section 3.1 we estimated the complete covariance matrix from scratch, using all selected steps from a *single generation*. We now take precisely the opposite viewpoint. We will repeatedly *update* the covariance matrix in the generation sequence using a *single selected step* only. First, this perspective will give a justification of the adaptation rule (14). Second, we will introduce the so-called evolution path that is finally used for a rank-one update of the covariance matrix.

3.3.1 A Different Viewpoint

We consider a specific method to produce n -dimensional normal distributions with zero mean. Let the vectors $\mathbf{y}_1, \dots, \mathbf{y}_{g_0} \in \mathbb{R}^n$, $g_0 \geq n$, span \mathbb{R}^n and let $\mathcal{N}(0, 1)$ denote independent $(0, 1)$ -normally distributed random numbers, then

$$\mathcal{N}(0, 1) \mathbf{y}_1 + \dots + \mathcal{N}(0, 1) \mathbf{y}_{g_0} \sim \mathcal{N}\left(\mathbf{0}, \sum_{i=1}^{g_0} \mathbf{y}_i \mathbf{y}_i^T\right) \quad (19)$$

is a normally distributed random vector with zero mean and covariance matrix $\sum_{i=1}^{g_0} \mathbf{y}_i \mathbf{y}_i^T$. The random vector (19) is generated by adding “line-distributions” $\mathcal{N}(0, 1) \mathbf{y}_i$. The singular distribution $\mathcal{N}(0, 1) \mathbf{y}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{y}_i \mathbf{y}_i^T)$ generates the vector \mathbf{y}_i with maximum likelihood considering all normal distributions with zero mean.

The line distribution that generates a vector \mathbf{y} with the maximum likelihood must “live” on a line that includes \mathbf{y} , and therefore the distribution must obey $\mathcal{N}(0, 1) \sigma \mathbf{y} \sim \mathcal{N}(0, \sigma^2 \mathbf{y} \mathbf{y}^T)$. Any other line distribution with zero mean cannot generate \mathbf{y} at all. Choosing σ reduces to choosing the maximum likelihood of $\|\mathbf{y}\|$ for the one-dimensional gaussian $\mathcal{N}(0, \sigma^2 \|\mathbf{y}\|^2)$, which is $\sigma = 1$.

¹³This can be shown more easily, because $(1 - c_{\text{cov}})^g = \exp \ln(1 - c_{\text{cov}})^g = \exp(g \ln(1 - c_{\text{cov}})) \approx \exp(-g c_{\text{cov}})$ for small c_{cov} , and for $g \approx 1/c_{\text{cov}}$ we get immediately $(1 - c_{\text{cov}})^g \approx \exp(-1)$.

¹⁴We use the sphere model $f_{\text{sphere}}(\mathbf{x}) = \sum_i x_i^2$ to empirically find a good setting for the parameter c_{cov} , dependent on n and μ_{eff} . The found setting was applicable to any non-noisy objective function we tried so far.

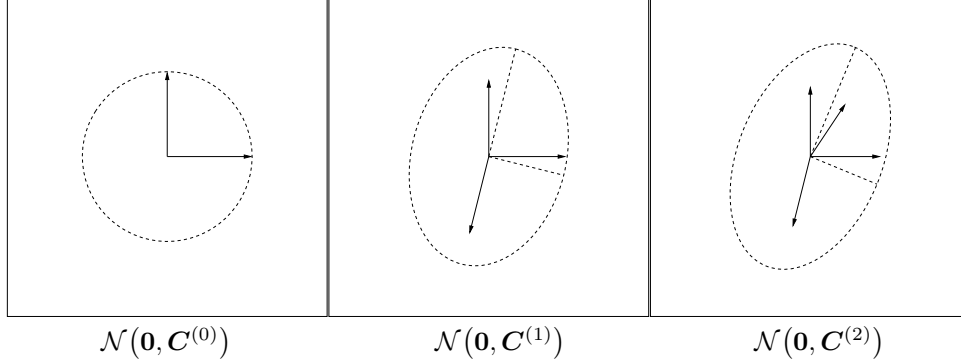


Figure 4: Change of the distribution according to the covariance matrix update (20). Left: vectors e_1 and e_2 , and $C^{(0)} = \mathbf{I} = e_1 e_1^T + e_2 e_2^T$. Middle: vectors $0.91 e_1$, $0.91 e_2$, and $0.41 \mathbf{y}_1$ (the coefficients deduce from $c_{\text{cov}} = 0.17$), and $C^{(1)} = (1 - c_{\text{cov}}) \mathbf{I} + c_{\text{cov}} \mathbf{y}_1 \mathbf{y}_1^T$, where $\mathbf{y}_1 = \begin{pmatrix} -0.59 \\ -2.2 \end{pmatrix}$. The distribution ellipsoid is elongated into the direction of \mathbf{y}_1 , and therefore increases the likelihood of \mathbf{y}_1 . Right: $C^{(2)} = (1 - c_{\text{cov}}) C^{(1)} + c_{\text{cov}} \mathbf{y}_2 \mathbf{y}_2^T$, where $\mathbf{y}_2 = \begin{pmatrix} 0.97 \\ 1.5 \end{pmatrix}$.

The covariance matrix $\mathbf{y}\mathbf{y}^T$ has rank one, its only eigenvectors are $\mathbb{R}_{\setminus 0} \times \mathbf{y}$ with eigenvalue $\|\mathbf{y}\|^2$. Using equation (19), any normal distribution can be realized if \mathbf{y}_i are chosen appropriately. For example, (19) resembles (4) with $\mathbf{m} = \mathbf{0}$, using the orthogonal eigenvectors $\mathbf{y}_i = d_{i,i} \mathbf{b}_i$, for $i = 1, \dots, n$, where \mathbf{b}_i are the columns of \mathbf{B} . In general, the vectors \mathbf{y}_i need not to be eigenvectors of the covariance matrix (and usually are not).

Considering (19) and a slight simplification of (14), we try to gain insight into the adaptation rule for the covariance matrix. Let the sum in (14) consist of a single summand only (e.g. $\mu = 1$), and let $\mathbf{y}_{g+1} = \frac{\mathbf{x}_{i,\lambda}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}}$. Then, the rank-one update for the covariance matrix reads

$$C^{(g+1)} = (1 - c_{\text{cov}}) C^{(g)} + c_{\text{cov}} \mathbf{y}_{g+1} \mathbf{y}_{g+1}^T \quad (20)$$

The right summand is of rank one and adds the maximum likelihood term for \mathbf{y}_{g+1} into the covariance matrix $C^{(g)}$. Therefore the probability to generate \mathbf{y}_{g+1} in the next generation increases.

An example of the first two iteration steps of (20) is shown in **Figure 4**. The distribution $\mathcal{N}(\mathbf{0}, C^{(1)})$ tends to reproduce \mathbf{y}_1 with a larger probability than the initial distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$; the distribution $\mathcal{N}(\mathbf{0}, C^{(2)})$ tends to reproduce \mathbf{y}_2 with a larger probability than $\mathcal{N}(\mathbf{0}, C^{(1)})$, and so forth. When $\mathbf{y}_1, \dots, \mathbf{y}_g$ denote the formerly selected, favorable steps, $\mathcal{N}(\mathbf{0}, C^{(g)})$ tends to reproduce these steps. The process leads to an alignment of the search distribution $\mathcal{N}(\mathbf{0}, C^{(g)})$ to the distribution of the selected steps. If both distributions become alike, as under random selection, in expectation no further change of the covariance matrix takes place [6].

3.3.2 Cumulation: Utilizing the Evolution Path

We have used the selected steps, $\mathbf{y}_{i:\lambda}^{(g+1)} = (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)})/\sigma^{(g)}$, to update the covariance matrix in (14) and (20). Because $\mathbf{y}\mathbf{y}^\top = -\mathbf{y}(-\mathbf{y})^\top$, the sign of the steps is irrelevant for the update of the covariance matrix—that is, the sign information is not used for calculating $\mathbf{C}^{(g+1)}$. To exploit the sign information, the so-called *evolution path* is introduced [12, 14].

We call a sequence of successive steps, the strategy takes over a number of generations, an evolution path. An evolution path can be expressed by a sum of consecutive steps. This summation is referred to as *cumulation*. To construct an evolution path, the step-size σ is disregarded. For example, an evolution path of three steps of the distribution mean \mathbf{m} can be constructed by the sum

$$\frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} + \frac{\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}}{\sigma^{(g-1)}} + \frac{\mathbf{m}^{(g-1)} - \mathbf{m}^{(g-2)}}{\sigma^{(g-2)}} . \quad (21)$$

In practice, to construct the evolution path, $\mathbf{p}_c \in \mathbb{R}^n$, we use exponential smoothing as in (14), and start with $\mathbf{p}_c^{(0)} = \mathbf{0}$.¹⁵

$$\mathbf{p}_c^{(g+1)} = (1 - c_c)\mathbf{p}_c^{(g)} + \sqrt{c_c(2 - c_c)\mu_{\text{eff}}} \frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} \quad (22)$$

where

$\mathbf{p}_c^{(g)} \in \mathbb{R}^n$, evolution path at generation g .

$c_c \leq 1$. Again, $1/c_c$ is the backward time horizon of the evolution path \mathbf{p}_c that contains roughly 63% of the overall weight (compare derivation of (18)). A time horizon between \sqrt{n} and n is reasonable.

The factor $\sqrt{c_c(2 - c_c)\mu_{\text{eff}}}$ is a normalization constant for \mathbf{p}_c . For $c_c = 1$ and $\mu_{\text{eff}} = 1$, the factor reduces to one, and $\mathbf{p}_c^{(g+1)} = (\mathbf{x}_{1:\lambda}^{(g+1)} - \mathbf{m}^{(g)})/\sigma^{(g)}$.

The factor $\sqrt{c_c(2 - c_c)\mu_{\text{eff}}}$ is chosen, such that

$$\mathbf{p}_c^{(g+1)} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad (23)$$

if

$$\mathbf{p}_c^{(g)} \sim \frac{\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \text{for all } i = 1, \dots, \mu . \quad (24)$$

To derive (23) from (24) and (22) remark that

$$(1 - c_c)^2 + \sqrt{c_c(2 - c_c)}^2 = 1 \quad \text{and} \quad \sum_{i=1}^{\mu} w_i \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \sim \frac{1}{\sqrt{\mu_{\text{eff}}}} \mathcal{N}(\mathbf{0}, \mathbf{C}) . \quad (25)$$

The (rank-one) update of the covariance matrix $\mathbf{C}^{(g)}$ via the evolution path $\mathbf{p}_c^{(g+1)}$ reads [12]

$$\mathbf{C}^{(g+1)} = (1 - c_{\text{cov}})\mathbf{C}^{(g)} + c_{\text{cov}}\mathbf{p}_c^{(g+1)}\mathbf{p}_c^{(g+1)\top} . \quad (26)$$

¹⁵In the final algorithm (22) is still slightly modified, compare (42).

An empirically validated choice for the learning rate in (26) is $c_{\text{cov}} \approx 2/n^2$. For $c_c = 1$ and $\mu = 1$, Equations (26), (20), and (14) are identical.

Using the evolution path for the update of \mathbf{C} is a significant improvement of (14) for small μ_{eff} , because correlations between consecutive steps are exploited. The leading signs of steps, and the dependencies between consecutive steps play a significant role for the resulting evolution path $\mathbf{p}_c^{(g+1)}$. For $c_c \approx 3/n$ the number of function evaluations needed to adapt a nearly optimal covariance matrix on cigar-like objective functions becomes $\mathcal{O}(n)$.

As a last step, we combine (14) and (26).

3.4 Combining Rank- μ -Update and Cumulation

The final CMA update of the covariance matrix combines (14) and (26), where μ_{cov} determines their relative weighting.

$$\begin{aligned} \mathbf{C}^{(g+1)} = & (1 - c_{\text{cov}})\mathbf{C}^{(g)} + \frac{c_{\text{cov}}}{\mu_{\text{cov}}} \underbrace{\mathbf{p}_c^{(g+1)}\mathbf{p}_c^{(g+1)\text{T}}}_{\text{rank-one update}} + c_{\text{cov}} \left(1 - \frac{1}{\mu_{\text{cov}}}\right) \\ & \times \underbrace{\sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}^{(g+1)} \left(\mathbf{y}_{i:\lambda}^{(g+1)}\right)^{\text{T}}}_{\text{rank-}\mu \text{ update}} \end{aligned} \quad (27)$$

where

$\mu_{\text{cov}} \geq 1$. Choosing $\mu_{\text{cov}} = \mu_{\text{eff}}$ is most appropriate.

$c_{\text{cov}} \approx \min(\mu_{\text{cov}}, \mu_{\text{eff}}, n^2)/n^2$.

$\mathbf{y}_{i:\lambda}^{(g+1)} = (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)})/\sigma^{(g)}$.

Equation (27) reduces to (14) for $\mu_{\text{cov}} \rightarrow \infty$ and to (26) for $\mu_{\text{cov}} = 1$. The equation combines the advantages of (14) and (26). On the one hand, the information within the population of one generation is used efficiently by the rank- μ update. On the other hand, information of correlations between generations is exploited by using the evolution path for the rank-one update. The former is important in large populations, the latter is in particular important in small populations.

4 Step-Size Control

The covariance matrix adaptation, introduced in the last section, does not explicitly control the “overall scale” of the distribution, the step-size. The covariance matrix adaptation increases the scale only in one direction for each selected step, and it decreases the scale only implicitly by fading out old information via the factor $1 - c_{\text{cov}}$. Less informally, we can state two specific reasons to introduce a step-size control in addition to the adaptation rule (27) for $\mathbf{C}^{(g)}$.

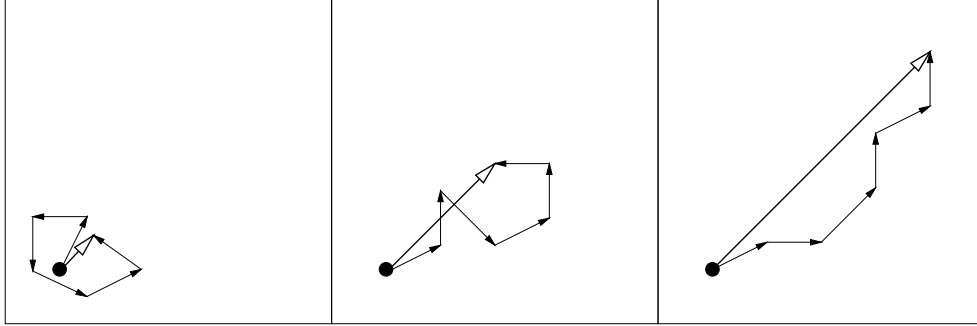


Figure 5: Three evolution paths of respectively six steps from different selection situations (idealized). The lengths of the *single* steps are all comparable. The length of the evolution paths (sum of steps) is remarkably different and is exploited for step-size control

1. The optimal overall step length cannot be well approximated by (27), in particular if μ_{eff} is chosen larger than one.

For example, on $f_{\text{sphere}}(\mathbf{x}) = \sum_{i=1}^n x_i^2$, the optimal step-size σ equals approximately $\mu \sqrt{f_{\text{sphere}}(\mathbf{x})}/n$, given $\mathbf{C}^{(g)} \approx \mathbf{I}$ and $\mu_{\text{eff}} = \mu \ll n$ [3, 18]. This dependency on μ cannot be realized by (14), and is also not well approximated by (27).

2. The largest reliable learning rate for the covariance matrix update in (27) is too slow to achieve competitive change rates for the overall step length.

To achieve optimal performance on f_{sphere} with an Evolution Strategy, the overall step length must decrease by a factor of approximately $\exp(0.202) \approx 1.22$ within n function evaluations, as can be derived from progress formulas [3, p. 229]. That is, the time horizon for the step length change must be proportional to n or shorter. From the learning rate c_{cov} in (27) follows that the adaptation is too slow to perform competitive on f_{sphere} whenever $\mu_{\text{eff}} \ll n$. This can be validated by simulations even for moderate dimensions, say, $n \geq 10$ and small μ_{eff} , say, $\mu_{\text{eff}} \leq 1 + \ln n$.

To control the step-size $\sigma^{(g)}$ we utilize an evolution path, *i.e.* a sum of successive steps (see page 15). The method can be applied independently of the covariance matrix update and is denoted as *cumulative path length control*, cumulative step-size control, or *cumulative step length adaptation (CSA)*. The length of an evolution path is exploited, based on the following reasoning (compare also Fig. 5).

- Whenever the evolution path is short, single steps cancel each other out (Fig. 5, left). Loosely speaking, they are anti-correlated. If steps annihilate each other, the step-size should be decreased.
- Whenever the evolution path is long, the single steps are pointing to similar directions (Fig. 5, right). Loosely speaking, they are correlated. Because the steps are similar, the same distance can be covered by fewer but longer steps into the same directions. In the

limit case, where consecutive steps have identical direction, they can be replaced by an enlarged single step. Consequently, the step-size should be increased.

- Subsuming, in the desired situation the steps are (approximately) perpendicular in expectation and therefore uncorrelated (Fig. 5, middle).

To decide whether the evolution path is “long” or “short”, we compare the length of the path with its *expected length under random selection*.¹⁶ Under random selection consecutive steps are independent and therefore uncorrelated (we just realized that “uncorrelated” steps are the desired situation). If selection biases the evolution path to be longer than expected, σ is increased, and, vice versa, if selection biases the evolution path to be shorter than expected, σ is decreased. In the ideal situation, selection does not bias the length of the evolution path and the length equals its expected length under random selection.

In practice, to construct the evolution path, \mathbf{p}_σ , the same techniques as in (22) are applied. In contrast to (22), a *conjugate* evolution path is constructed, because the expected length of the evolution path \mathbf{p}_c from (22) depends on its direction (compare (23)). Initialized with $\mathbf{p}_\sigma^{(0)} = \mathbf{0}$, the conjugate evolution path reads

$$\mathbf{p}_\sigma^{(g+1)} = (1 - c_\sigma)\mathbf{p}_\sigma^{(g)} + \sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}} \mathbf{C}^{(g)-\frac{1}{2}} \frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} \quad (28)$$

where

$\mathbf{p}_\sigma^{(g)} \in \mathbb{R}^n$ is the conjugate evolution path at generation g .

$c_\sigma < 1$. Again, $1/c_\sigma$ is the backward time horizon of the evolution path (compare (18)). For small μ_{eff} , a time horizon between \sqrt{n} and n is reasonable.

$\sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}}$ is a normalization constant, see (22).

$\mathbf{C}^{(g)-\frac{1}{2}} \stackrel{\text{def}}{=} \mathbf{B}^{(g)} \mathbf{D}^{(g)-1} \mathbf{B}^{(g)\text{T}}$, where $\mathbf{C}^{(g)} = \mathbf{B}^{(g)} (\mathbf{D}^{(g)})^2 \mathbf{B}^{(g)\text{T}}$ is an eigendecomposition of $\mathbf{C}^{(g)}$, where $\mathbf{B}^{(g)}$ is an orthonormal basis of eigenvectors, and the diagonal elements of the diagonal matrix $\mathbf{D}^{(g)}$ are square roots of the corresponding positive eigenvalues (cf. Sect. 0.1).

For $\mathbf{C}^{(g)} = \mathbf{I}$, we have $\mathbf{C}^{(g)-\frac{1}{2}} = \mathbf{I}$ and (28) replicates (22). The transformation $\mathbf{C}^{(g)-\frac{1}{2}}$ re-scales the step $\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}$ within the coordinate system given by $\mathbf{B}^{(g)}$.

The single factors of the transformation $\mathbf{C}^{(g)-\frac{1}{2}} = \mathbf{B}^{(g)} \mathbf{D}^{(g)-1} \mathbf{B}^{(g)\text{T}}$ can be explained as follows (from right to left):

$\mathbf{B}^{(g)\text{T}}$ rotates the space such that the columns of $\mathbf{B}^{(g)}$, *i.e.* the principle axes of the distribution $\mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)})$, rotate into the coordinate axes. Elements of the resulting vector relate to projections onto the corresponding eigenvectors.

$\mathbf{D}^{(g)-1}$ applies a (re-)scaling such that all axes become equally sized.

¹⁶Random selection means that the index $i : \lambda$ (compare (6)) is independent of the value of $\mathbf{x}_{i:\lambda}^{(g+1)}$ for all $i = 1, \dots, \lambda$, e.g. $i : \lambda = i$.

$B^{(g)}$ rotates the result back into the original coordinate system. This last transformation ensures that the principal axes of the distribution are not rotated by the overall transformation and directions of consecutive steps are comparable.

Consequently, the transformation $C^{(g)-\frac{1}{2}}$ makes the expected length of $\mathbf{p}_\sigma^{(g+1)}$ independent of its direction, and for any sequence of realized covariance matrices $C_{g=0,1,2,\dots}^{(g)}$ we have under random selection $\mathbf{p}_\sigma^{(g+1)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, given $\mathbf{p}_\sigma^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ [6].

To update $\sigma^{(g)}$, we “compare” $\|\mathbf{p}_\sigma^{(g+1)}\|$ with its expected length $E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$, that is

$$\ln \sigma^{(g+1)} = \ln \sigma^{(g)} + \frac{c_\sigma}{d_\sigma E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} \left(\|\mathbf{p}_\sigma^{(g+1)}\| - E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\| \right), \quad (29)$$

where

$d_\sigma \approx 1$, damping parameter, scales the change magnitude of $\ln \sigma^{(g)}$. The factor $c_\sigma/d_\sigma/E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$ is based on in-depth investigations of the algorithm [6].

$E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\| = \sqrt{2}\Gamma(\frac{n+1}{2})/\Gamma(\frac{n}{2}) \approx \sqrt{n} + \mathcal{O}(1/n)$, expectation of the Euclidean norm of a $\mathcal{N}(\mathbf{0}, \mathbf{I})$ distributed random vector.

For $\|\mathbf{p}_\sigma^{(g+1)}\| = E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$ the second summand in (29) is zero, and $\sigma^{(g)}$ is unchanged, while $\sigma^{(g)}$ is increased for $\|\mathbf{p}_\sigma^{(g+1)}\| > E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$, and $\sigma^{(g)}$ is decreased for $\|\mathbf{p}_\sigma^{(g+1)}\| < E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$.

Alternatively, we might use the squared norm $\|\mathbf{p}_\sigma^{(g+1)}\|^2$ in (29) and compare with its expected value n [2]. In this case (29) would read

$$\begin{aligned} \ln \sigma^{(g+1)} &= \ln \sigma^{(g)} + \frac{c_\sigma}{d_\sigma 2n} \left(\|\mathbf{p}_\sigma^{(g+1)}\|^2 - n \right) \\ &= \ln \sigma^{(g)} + \frac{c_\sigma}{2d_\sigma} \left(\frac{\|\mathbf{p}_\sigma^{(g+1)}\|^2}{n} - 1 \right). \end{aligned} \quad (30)$$

This update will presumably lead to faster step-size increments and slower step-size decrements.

The step-size change is unbiased on the log scale, because $E[\ln \sigma^{(g+1)} | \sigma^{(g)}] = \ln \sigma^{(g)}$ for $\mathbf{p}_\sigma^{(g+1)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The role of unbiasedness is discussed in Sect. 5. Equations (28) and (29) cause successive steps of the distribution mean $\mathbf{m}^{(g)}$ to be approximately $C^{(g)-1}$ -conjugate.

In order to show that successive steps are approximately $C^{(g)-1}$ -conjugate first we remark that (28) and (29) adapt σ such that the length of $\mathbf{p}_\sigma^{(g+1)}$ equals approximately $E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$. Starting from $(E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|)^2 \approx \|\mathbf{p}_\sigma^{(g+1)}\|^2 = \mathbf{p}_\sigma^{(g+1)\top} \mathbf{p}_\sigma^{(g+1)} = \text{RHS}^\top \text{RHS}$ of (28) and assuming that the expected squared length of $C^{(g)-\frac{1}{2}}(\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)})$ is unchanged by selection (unlike its direction) we get

$$\mathbf{p}_\sigma^{(g)\top} C^{(g)-\frac{1}{2}}(\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}) \approx 0, \quad (31)$$

and

$$\left(\mathbf{C}^{(g)\frac{1}{2}}\mathbf{p}_\sigma^{(g)}\right)^T \mathbf{C}^{(g)-1} \left(\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}\right) \approx 0 . \quad (32)$$

Given $1/c_{cov} \gg 1$ and (31) we assume also $\mathbf{p}_\sigma^{(g-1)T} \mathbf{C}^{(g)-\frac{1}{2}} \left(\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}\right) \approx 0$ and derive

$$\left(\mathbf{m}^{(g)} - \mathbf{m}^{(g-1)}\right)^T \mathbf{C}^{(g)-1} \left(\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}\right) \approx 0 . \quad (33)$$

That is, the steps taken by the distribution mean become approximately $\mathbf{C}^{(g)-1}$ -conjugate.

Because $\sigma^{(g)} > 0$, (29) is equivalent to

$$\sigma^{(g+1)} = \sigma^{(g)} \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma^{(g+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right) \quad (34)$$

The length of the evolution path is an intuitive and empirically well validated goodness measure for the overall step length. For $\mu_{\text{eff}} > 1$ it is the best measure to our knowledge. Nevertheless, it fails to adapt nearly optimal step-sizes on very noisy objective functions [4].

5 Discussion

The CMA-ES is an attractive option for non-linear optimization, if ‘‘classical’’ search methods, e.g. quasi-Newton methods (BFGS) and/or conjugate gradient methods, fail due to a non-convex or rugged search landscape (e.g. sharp bends, discontinuities, outliers, noise, and local optima). Learning the covariance matrix in the CMA-ES is analogous to learning the inverse Hessian matrix in a quasi-Newton method. In the end, any convex-quadratic (ellipsoid) objective function is transformed into the spherical function f_{sphere} . This can improve the performance on ill-conditioned and/or non-separable problems by orders of magnitude.

The CMA-ES overcomes typical problems that are often associated with evolutionary algorithms.

1. Poor performance on badly scaled and/or highly non-separable objective functions. Equation (27) adapts the search distribution to badly scaled and non-separable problems.
2. The inherent need to use large population sizes. A typical, however intricate to diagnose reason for the failure of population based search algorithms is the degeneration of the population into a subspace. This is usually prevented by non-adaptive components in the algorithm and/or by a large population size (considerably larger than the problem dimension). In the CMA-ES, the population size can be freely chosen, because the learning rate c_{cov} in (27) prevents the degeneration even for small population sizes, e.g. $\lambda = 9$. Small population sizes usually lead to faster convergence, large population sizes help to avoid local optima.
3. Premature convergence of the population. Step-size control in (34) prevents the population to converge prematurely. It does not prevent the search to end up in a local optimum.

Therefore, the CMA-ES is highly competitive on a considerable number of test functions [6, 9, 11, 13, 14] and was successfully applied to many real world problems.¹⁷ Finally we discuss a few **basic design principles** that were applied in the previous sections.

Change rates We refer to a change rate as the expected parameter change *per sampled search point*, given a certain selection situation. To achieve competitive performance on a wide range of objective functions, the possible change rates of the adaptive parameters need to be adjusted carefully. The CMA-ES separately controls change rates for the mean value of the distribution, \mathbf{m} , the covariance matrix, \mathbf{C} , and the step-size, σ .

- The change rate for the mean value \mathbf{m} , given a fixed sample distribution, is determined by the parent number and the recombination weights. The larger μ_{eff} , the smaller is the possible change rate of \mathbf{m} . Similar holds for most evolutionary algorithms.
- The change rate of the covariance matrix \mathbf{C} is explicitly controlled by the learning rate c_{cov} and therefore detached from parent number and population size. The learning rate reflects the model complexity. In evolutionary algorithms, the explicit control of change rates of covariances, independent from the population size, is a rare feature.
- The change rate of the step-size σ is explicitly controlled by the damping parameter d_σ and is in particular independent from the change rate of \mathbf{C} . The time constant $1/c_\sigma \leq n$ ensures a sufficiently fast change of the overall step length in particular with small population sizes.

Invariance Invariance properties of a search algorithm denote identical behavior on a set, or a class of objective functions. Invariance is an important property of the CMA-ES.¹⁸ Translation invariance should be taken for granted in continuous domain optimization. Translation invariance means that the search behavior on the function $\mathbf{x} \mapsto f(\mathbf{x} + \mathbf{a})$, $\mathbf{x}^{(0)} = \mathbf{b} - \mathbf{a}$, is independent of $\mathbf{a} \in \mathbb{R}^n$. Further invariances, e.g. to certain linear transformations of the search space, are highly desirable: they imply uniform performance on classes of functions and therefore allow for generalization of empirical results. The CMA-ES exhibits the following invariances.

- Invariance to order preserving (*i.e.* strictly monotonic) transformations of the objective function value. The algorithm only depends on *the ranking* of function values.
- Invariance to angle preserving (rigid) transformations of the search space (rotation, reflection, and translation) if the initial search point is transformed accordingly.
- Scale invariance if the initial scaling, e.g. $\sigma^{(0)}$, and the initial search point, $\mathbf{m}^{(0)}$, are chosen accordingly.
- Invariance to a scaling of variables if the initial diagonal covariance matrix $\mathbf{C}^{(0)}$, and the initial search point, $\mathbf{m}^{(0)}$, are chosen accordingly.

¹⁷For a list of published applications see <http://www.lri.fr/~hansen/cmaapplications.pdf>

¹⁸Special acknowledgments to Iván Santibáñez-Koref for pointing this out to me.

- Invariance to any invertible linear transformation of the search space, \mathbf{A} , if the initial covariance matrix $\mathbf{C}^{(0)} = \mathbf{A}^{-1} (\mathbf{A}^{-1})^T$, and if the initial search point, $\mathbf{m}^{(0)}$, are transformed accordingly.

Invariance should be a fundamental design criterion for any search algorithm. Together with the ability to efficiently adapt the invariance governing parameters, invariance is a key to competitive performance.

Stationarity An important design criterion for a stochastic search procedure is *unbiasedness* of variations of object and strategy parameters [5, 14]. Consider random selection, e.g. the objective function $f(\mathbf{x}) = \text{rand}$ to be independent of \mathbf{x} . The population mean is unbiased if its expected value remains unchanged in the next generation, that is $E[\mathbf{m}^{(g+1)} | \mathbf{m}^{(g)}] = \mathbf{m}^{(g)}$. For the population mean, stationarity under random selection is a rather intuitive concept. In the CMA-ES, stationarity is respected for all parameters in the basic equation (5). The distribution mean \mathbf{m} , the covariance matrix \mathbf{C} , and $\ln \sigma$ are unbiased. Unbiasedness of $\ln \sigma$ does not imply that σ is unbiased. Actually, under random selection, $E[\sigma^{(g+1)} | \sigma^{(g)}] > \sigma^{(g)}$, compare (29).¹⁹

For distribution variances (or step-sizes) a bias toward increase or decrease entails the danger of divergence or premature convergence, respectively, whenever the selection pressure is low. Nevertheless, on noisy problems a properly controlled bias towards increase might be appropriate even on the log scale.

Acknowledgments

The author wishes to gratefully thank Anne Auger, Christian Igel, Stefan Kern, and Fabrice Marchal for the many valuable comments on the manuscript.

References

- [1] Auger A, Hansen N. A restart CMA evolution strategy with increasing population size. In *Proceedings of the IEEE Congress on Evolutionary Computation*, 2005.
- [2] Arnold DV, Beyer HG. Performance analysis of evolutionary optimization with cumulative step length adaptation. *IEEE Transactions on Automatic Control*, 49(4):617–622, 2004.
- [3] Beyer HG. *The Theory of Evolution Strategies*. Springer, Berlin, 2001.
- [4] Beyer HG, Arnold DV. Qualms regarding the optimality of cumulative path length control in CSA/CMA-evolution strategies. *Evolutionary Computation*, 11(1):19–28, 2003.
- [5] Beyer HG, Deb K. On self-adaptive features in real-parameter evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 5(3):250–270, 2001.

¹⁹Alternatively, if (34) would be designed to be unbiased for $\sigma^{(g+1)}$, this would imply that $E[\ln \sigma^{(g+1)} | \sigma^{(g)}] < \ln \sigma^{(g)}$, to our opinion a less desirable variant.

- [6] Hansen N. *Verallgemeinerte individuelle Schrittweitenregelung in der Evolutionsstrategie*. Mensch und Buch Verlag, Berlin, 1998.
- [7] Hansen N. Invariance, self-adaptation and correlated mutations in evolution strategies. In Schoenauer M, Deb K, Rudolph G, Yao X, Lutton E, Merelo JJ, Schwefel HP, editors, *Parallel Problem Solving from Nature - PPSN VI*, pages 355–364. Springer, 2000.
- [8] Hansen N. The CMA evolution strategy: a comparing review. In Lozano JA, Larranaga P, Inza I, and Bengoetxea E, editors, *Towards a new evolutionary computation. Advances on estimation of distribution algorithms*, pages 75–102. Springer, 2006.
- [9] Hansen N, Kern S. Evaluating the CMA evolution strategy on multimodal test functions. In Xin Yao et al., editors, *Parallel Problem Solving from Nature - PPSN VIII*, pages 282–291. Springer, 2004.
- [10] Hansen N, Niederberger SPN, Guzzella L, Koumoutsakos P. A method for handling uncertainty in evolutionary optimization with an application to feedback control of combustion. *IEEE Transactions on Evolutionary Computation*, 2009.
- [11] Hansen N, Müller SD, Koumoutsakos P. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). *Evolutionary Computation*, 11(1):1–18, 2003.
- [12] Hansen N, Ostermeier A. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. In *Proceedings of the 1996 IEEE Conference on Evolutionary Computation (ICEC '96)*, pages 312–317, 1996.
- [13] Hansen N, Ostermeier A. Convergence properties of evolution strategies with the derandomized covariance matrix adaptation: The $(\mu/\mu_I, \lambda)$ -CMA-ES. In *Proceedings of the 5th European Congress on Intelligent Techniques and Soft Computing*, pages 650–654, 1997.
- [14] Hansen N, Ostermeier A. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.
- [15] Kern S, Müller SD, Hansen N, Büche D, Ocenasek J, Koumoutsakos P. Learning probability distributions in continuous evolutionary algorithms – a comparative review. *Natural Computing*, 3:77–112, 2004.
- [16] Larrañaga P. A review on estimation of distribution algorithms. In P. Larrañaga and J. A. Lozano, editors, *Estimation of Distribution Algorithms*, pages 80–90. Kluwer Academic Publishers, 2002.
- [17] Larrañaga P, Lozano JA, Bengoetxea E. Estimation of distribution algorithms based on multivariate normal and Gaussian networks. Technical report, Dept. of Computer Science and Artificial Intelligence, University of the Basque Country, 2001. KZAA-IK-1-01.
- [18] Rechenberg I. *Evolutionsstrategie '94*. Frommann-Holzboog, Stuttgart, Germany, 1994.
- [19] Rubenstein RY, Kroese DP. The Cross-Entropy Method: a unified approach to combinatorial optimization, Monte-Carlo simulation, and machine learning. Springer, 2004.

A Algorithm Summary: The $(\mu/\mu_W, \lambda)$ -CMA-ES

Figure 6 outlines the complete algorithm, summarizing (5), (6), (22), (27), (28), and (34). Used symbols are:

$\mathbf{y}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$, for $k = 1, \dots, \lambda$, are realizations from a multivariate normal distribution with zero mean and covariance matrix \mathbf{C} .

\mathbf{B}, \mathbf{D} result from an eigendecomposition of the covariance matrix \mathbf{C} with $\mathbf{C} = \mathbf{B}\mathbf{D}^2\mathbf{B}^T = \mathbf{B}\mathbf{D}\mathbf{D}\mathbf{B}^T$ (cf. Sect. 0.1). Columns of \mathbf{B} are an orthonormal basis of eigenvectors. Diagonal elements of the diagonal matrix \mathbf{D} are square roots of the corresponding positive eigenvalues. While (36) can certainly be implemented using a Cholesky decomposition of \mathbf{C} the eigendecomposition is needed to compute (40) correctly.

$\mathbf{x}_k \in \mathbb{R}^n$, for $k = 1, \dots, \lambda$. Sample of λ search points.

$\langle \mathbf{y} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}$, step of the distribution mean disregarding step-size σ .

$\mathbf{y}_{i:\lambda} = (\mathbf{x}_{i:\lambda} - \mathbf{m})/\sigma$, see $\mathbf{x}_{i:\lambda}$ below.

$\mathbf{x}_{i:\lambda} \in \mathbb{R}^n$, i -th best point out of $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$ from (37). The index $i : \lambda$ denotes the index of the i -th ranked point, that is $f(\mathbf{x}_{1:\lambda}) \leq f(\mathbf{x}_{2:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.

$\mu_{\text{eff}} = (\sum_{i=1}^{\mu} w_i^2)^{-1}$ is the variance effective selection mass, see (8). It holds $1 \leq \mu_{\text{eff}} \leq \mu$.

$\mathbf{C}^{-\frac{1}{2}} \stackrel{\text{def}}{=} \mathbf{B}\mathbf{D}^{-1}\mathbf{B}^T$, see \mathbf{B}, \mathbf{D} above. The matrix \mathbf{D} can be inverted by inverting its diagonal elements. From the definitions we find that $\mathbf{C}^{-\frac{1}{2}} \langle \mathbf{y} \rangle_w = \mathbf{B} \langle \mathbf{z} \rangle_w$ with $\langle \mathbf{z} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$.

$\mathbb{E} \|\mathcal{N}(\mathbf{0}, \mathbf{I})\| = \sqrt{2} \Gamma(\frac{n+1}{2}) / \Gamma(\frac{n}{2}) \approx \sqrt{n} (1 - \frac{1}{4n} + \frac{1}{21n^2})$.

$$h_\sigma = \begin{cases} 1 & \text{if } \frac{\|\mathbf{p}_\sigma\|}{\sqrt{1-(1-c_\sigma)^{2(g+1)}}} < (1.4 + \frac{2}{n+1}) \mathbb{E} \|\mathcal{N}(\mathbf{0}, \mathbf{I})\| \\ 0 & \text{otherwise} \end{cases}$$
, where g is the generation

number. The Heaviside function h_σ stalls the update of \mathbf{p}_c in (42) if $\|\mathbf{p}_\sigma\|$ is large. This prevents a too fast increase of axes of \mathbf{C} in a linear surrounding, *i.e.* when the step-size is far too small. This is useful when the initial step-size is chosen far too small or when the objective function changes in time.

$\delta(h_\sigma) = (1 - h_\sigma)c_c(2 - c_c) \leq 1$ is of minor relevance. In the (unusual) case of $h_\sigma = 0$, it substitutes for the second summand from (42) in (43).

Default Parameters The (external) strategy parameters are $\lambda, \mu, w_{i=1\dots\mu}, c_\sigma, d_\sigma, c_c, \mu_{\text{cov}}$, and c_{cov} . Default strategy parameters values are given in Table 1. An in-depth discussion of most parameters is given in [14]. The default parameters of (46)-(49) are in particular chosen to be a robust setting and therefore, to our experience, applicable to a wide range of functions to be optimized. We do not recommend to change this setting. In contrast, the population size

Set parameters

Set parameters $\lambda, \mu, w_{i=1\dots\mu}, c_\sigma, d_\sigma, c_c, \mu_{\text{cov}},$ and c_{cov} to their default values according to Table 1.

Initialization

Set evolution paths $\mathbf{p}_\sigma = \mathbf{0}, \mathbf{p}_c = \mathbf{0}$, covariance matrix $\mathbf{C} = \mathbf{I}$, and $g = 0$.

Choose distribution mean $\mathbf{m} \in \mathbb{R}^n$ and step-size $\sigma \in \mathbb{R}_+$ problem dependent.¹

Until termination criterion met, $g \leftarrow g + 1$

Sample new population of search points, for $k = 1, \dots, \lambda$

$$\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (35)$$

$$\mathbf{y}_k = \mathbf{B}\mathbf{D}\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad (36)$$

$$\mathbf{x}_k = \mathbf{m} + \sigma\mathbf{y}_k \sim \mathcal{N}(\mathbf{m}, \sigma^2\mathbf{C}) \quad (37)$$

Selection and recombination

$$\langle \mathbf{y} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \quad \text{where } \sum_{i=1}^{\mu} w_i = 1, w_i > 0 \quad (38)$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{y} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} \quad (39)$$

Step-size control

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma)\mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}} \mathbf{C}^{-\frac{1}{2}} \langle \mathbf{y} \rangle_w \quad (40)$$

$$\sigma \leftarrow \sigma \times \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right) \quad (41)$$

Covariance matrix adaptation

$$\mathbf{p}_c \leftarrow (1 - c_c)\mathbf{p}_c + h_\sigma \sqrt{c_c(2 - c_c)\mu_{\text{eff}}} \langle \mathbf{y} \rangle_w \quad (42)$$

$$\begin{aligned} \mathbf{C} \leftarrow & (1 - c_{\text{cov}})\mathbf{C} + \frac{c_{\text{cov}}}{\mu_{\text{cov}}} (\mathbf{p}_c \mathbf{p}_c^\top + \delta(h_\sigma)\mathbf{C}) \\ & + c_{\text{cov}} \left(1 - \frac{1}{\mu_{\text{cov}}}\right) \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^\top \end{aligned} \quad (43)$$

¹The optimum should presumably be within the initial cube $\mathbf{m} \pm 3\sigma(1, \dots, 1)^\top$. If the optimum is expected to be in the initial search interval $[a, b]^n$ we may choose the initial search point, \mathbf{m} , uniformly randomly in $[a, b]^n$, and $\sigma = 0.3(b - a)$. Different search intervals Δs_i for different variables can be reflected by a different initialization of \mathbf{C} , in that the diagonal elements of \mathbf{C} obey $c_{ii} = (\Delta s_i)^2$. Remark that the Δs_i should not disagree by several orders of magnitude. Otherwise a scaling of the variables should be applied.

Figure 6: The $(\mu/\mu_W, \lambda)$ -CMA Evolution Strategy. Symbols: see text

Table 1: Default Strategy Parameters, where $\mu_{\text{eff}} = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \geq 1$ and $\sum_{i=1}^{\mu} w_i = 1$, taken from [9], where $w'_i = \ln(\mu + 1) - \ln i$ was used

Selection and Recombination:	
$\lambda = 4 + \lfloor 3 \ln n \rfloor, \quad \mu = \lceil \mu' \rceil, \quad \mu' = \frac{\lambda}{2}$ (44)	
$w_i = \frac{w'_i}{\sum_{j=1}^{\mu} w'_j}, \quad w'_i = \ln(\mu' + 0.5) - \ln i \quad \text{for } i = 1, \dots, \mu,$ (45)	
Step-size control:	
$c_{\sigma} = \frac{\mu_{\text{eff}} + 2}{n + \mu_{\text{eff}} + 3}, \quad d_{\sigma} = 1 + 2 \max\left(0, \sqrt{\frac{\mu_{\text{eff}} - 1}{n + 1}} - 1\right) + c_{\sigma}$ (46)	
Covariance matrix adaptation:	
$c_c = \frac{4}{n + 4}$ (47)	
$\mu_{\text{cov}} = \mu_{\text{eff}}$ (48)	
$c_{\text{cov}} = \frac{1}{\mu_{\text{cov}}} \frac{2}{(n + \sqrt{2})^2} + \left(1 - \frac{1}{\mu_{\text{cov}}}\right) \min\left(1, \frac{2\mu_{\text{cov}} - 1}{(n + 2)^2 + \mu_{\text{cov}}}\right)$ (49)	

λ in (44) can be increased.²⁰ If the λ -dependent default values for μ and w_i are used, the population size λ has a significant influence on the global search performance [9]. Increasing λ usually improves the global search capability and the robustness of the CMA-ES, at the price of a reduced convergence speed. The convergence speed decrease at most linearly with λ . Independent restarts with increasing population size [1], automated or manually conducted, are an appropriate policy to perform well on most problems.

A.1 Implementational Concerns

We discuss a few implementational concerns.

Multivariate normal distribution: Let the vector $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ have independent, $(0, 1)$ -normally distributed components that can easily be sampled on a computer. To generate a random vector $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ for (36), we set $\mathbf{y} = \mathbf{B}\mathbf{D}\mathbf{z}$ (see above symbol descriptions of \mathbf{B} and \mathbf{D} and Sects. 0.1 and 0.2, and compare lines 52–53 and 83–84 in the source code below). Given $\mathbf{y}_k = \mathbf{B}\mathbf{D}\mathbf{z}_k$ and $\mathbf{C}^{-\frac{1}{2}} = \mathbf{B}\mathbf{D}^{-1}\mathbf{B}^T$ we have $\mathbf{C}^{-\frac{1}{2}} \langle \mathbf{y} \rangle_{\text{w}} = \mathbf{B} \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$ (compare (40) and lines 61 and 64 in the source code below).

Strategy internal numerical effort: In practice, the re-calculation of \mathbf{B} and \mathbf{D} needs to be done not until $\max(1, \lfloor 1/(10nc_{\text{cov}}) \rfloor)$ generations. For reasonable c_{cov} values, this

²⁰Decreasing λ is not recommended. Too small values regularly have strong adverse effects on the performance.

reduces the numerical effort due to the eigendecomposition from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$ per generated search point, that is the effort of a matrix vector multiplication.

On a Pentium 4, 2.5 GHz processor the overall strategy internal time consumption is roughly $3 \times (n + 5)^2 \times 10^{-8}$ seconds per function evaluation [15].

Remark that it is not sufficient to compute a Cholesky decomposition of \mathbf{C} , because then (40) cannot be computed correctly.

Termination criteria: In general, the algorithm should be stopped whenever it becomes a waste of CPU-time to continue, and it would be better to restart (eventually with increased population size [1]) or to reconsidering the encoding and/or objective function formulation. We recommend the following termination criteria [1] that are mostly related to numerical stability:

- `NoEffectAxis`: stop if adding a 0.1-standard deviation vector in any principal axis direction of \mathbf{C} does not change \mathbf{m} .²¹
- `NoEffectCoord`: stop if adding 0.2-standard deviations in any single coordinate does not change \mathbf{m} (i.e. m_i equals $m_i + 0.2 \sigma c_{i,i}$ for any i).
- `EqualFunValues`: stop if the range of the best objective function values of the last $10 + \lceil 30n/\lambda \rceil$ generations is zero.
- `ConditionCov`: stop if the condition number of the covariance matrix exceeds 10^{14} .
- `TolXUp`: stop if $\sigma \times \max(\text{diag}(\mathbf{D}))$ increased by more than 10^4 . This usually indicates a far too small initial σ , or divergent behavior.

Two other useful termination criteria should be considered problem dependent:

- `TolFun`: stop if the range of the best objective function values of the last $10 + \lceil 30n/\lambda \rceil$ generations and all function values of the recent generation is below `TolFun`. Choosing `TolFun` depends on the problem, while 10^{-12} is a conservative first guess.
- `TolX`: stop if the standard deviation of the normal distribution is smaller than in all coordinates and σp_c is smaller than `TolX` in all components. By default we set `TolX` to 10^{-12} times the initial σ .

Flat fitness: In the case of equal function values for several individuals in the population, it is feasible to increase the step-size (see lines 92–96 in the source code below). This method can interfere with the termination criterion `TolFun`. In practice, observation of a flat fitness should be rather a termination criterion and consequently lead to a re-consideration of the objective function formulation.

Boundaries and Constraints: The handling of boundaries and constraints is problem dependent.

²¹More formally, we terminate if \mathbf{m} equals to $\mathbf{m} + 0.1 \sigma d_{ii} \mathbf{b}_i$, where $i = (g \bmod n) + 1$, and d_{ii}^2 and \mathbf{b}_i are respectively the i -th eigenvalue and eigenvector of \mathbf{C} , with $\|\mathbf{b}_i\| = 1$.

1. If the optimal solution is *not close to the infeasible domain*, a simple and sufficient way to handle any type of boundaries and constraints is re-sampling any infeasible solution \mathbf{x} until it become feasible. One might also choose the set the fitness as

$$f_{\text{fitness}}(\mathbf{x}) = f_{\text{max}} + \alpha \|\mathbf{x} - \mathbf{x}_{\text{feasible}}\| , \quad (50)$$

where f_{max} is larger than the worst fitness in the feasible population (in case of minization) and $\mathbf{x}_{\text{feasible}}$ is a constant feasible point, preferably in the middle of the feasible domain.

2. Otherwise, if a repair mechanism is available (as for box-constraints) we evaluate the objective function on a repaired search point, $\mathbf{x}_{\text{repaired}}$, and add a penalty depending on the distance to the repaired solution.

$$f_{\text{fitness}}(\mathbf{x}) = f(\mathbf{x}_{\text{repaired}}) + \alpha \|\mathbf{x} - \mathbf{x}_{\text{repaired}}\|^2 . \quad (51)$$

The repaired solution is disregarded afterwards.

In case of box-boundaries, $\mathbf{x}_{\text{repaired}}$ is set to the feasible solution with the smallest distance $\|\mathbf{x} - \mathbf{x}_{\text{repaired}}\|$. In other words, components that are infeasible in \mathbf{x} are set to the (closest) boundary value in $\mathbf{x}_{\text{repaired}}$. A similar boundary handling with a component-wise adaptive α is described in [10].

Otherwise, when no repair mechanism is available, the fitness of the infeasible search point \mathbf{x} might similarly compute to

$$f_{\text{fitness}}(\mathbf{x}) = f_{\text{offset}} + \alpha \sum_i \mathbb{1}_{c_i > 0} \times c_i(\mathbf{x})^2 \quad (52)$$

where, w.l.o.g., the (non-linear) constraints $c_i : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto c_i(\mathbf{x})$ are satisfied for $c_i(\mathbf{x}) \leq 0$, and the indicator function $\mathbb{1}_{c_i > 0}$ equals to one for $c_i(\mathbf{x}) > 0$, zero otherwise, and $f_{\text{offset}} = \text{median}_k f(\mathbf{x}_k)$ equals to the median function value of the feasible points in the same generation. This approach has not yet been experimentally evaluated by the author.

In either case of (51) and (52), α should be chosen such that the differences in f and the differences in the second summand have a similar magnitude.

B MATLAB Source Code

```

1 function xmin=purecmaes
2 % CMA-ES: Evolution Strategy with Covariance Matrix Adaptation for
3 % nonlinear function minimization.
4 %
5 % This code is an excerpt from cmaes.m and implements the key parts
6 % of the algorithm. It is intended to be used for READING and
7 % UNDERSTANDING the basic flow and all details of the CMA *algorithm*.
8 % Computational efficiency is sometimes disregarded.
9
10 % ----- Initialization -----
11
12 % User defined input parameters (need to be edited)
13 strfitnessfct = 'felli'; % name of objective/fitness function
14 N = 10; % number of objective variables/problem dimension
15 xmean = rand(N,1); % objective variables initial point
16 sigma = 0.5; % coordinate wise standard deviation (step-size)
17 stopfitness = 1e-10; % stop if fitness < stopfitness (minimization)
18 stopeval = 1e3*N^2; % stop after stopeval number of function evaluations
19
20 % Strategy parameter setting: Selection
21 lambda = 4+floor(3*log(N)); % population size, offspring number
22 mup = (lambda-1)/2; % lambda=12; mu=3; weights = ones(mu,1); would be (3_I,12)-ES
23 mu = ceil(mup); % number of parents/points for recombination
24 weights = log(mup+1)-log(1:mu)'; % muXone recombination weights
25 weights = weights/sum(weights); % normalize recombination weights array
26 mueff=sum(weights)^2/sum(weights.^2); % variance-effective size of mu
27
28 % Strategy parameter setting: Adaptation
29 cc = 4/(N+4); % time constant for cumulation for covariance matrix
30 cs = (mueff+2)/(N+mueff+3); % t-const for cumulation for sigma control
31 mucov = mueff; % size of mu used for calculating learning rate ccov
32 ccov = (1/mucov) * 2/(N+1.4)^2 + (1-1/mucov) * ... % learning rate for
33 (2*mueff-1)/((N+2)^2+2*mueff); % covariance matrix
34 damp = 1 + 2*max(0, sqrt((mueff-1)/(N+1))-1) + cs; % damping for sigma
35
36 % Initialize dynamic (internal) strategy parameters and constants
37 pc = zeros(N,1); ps = zeros(N,1); % evolution paths for C and sigma
38 B = eye(N); % B defines the coordinate system
39 D = eye(N); % diagonal matrix D defines the scaling
40 C = B*D*(B*D)'; % covariance matrix
41 eigeneval = 0; % B and D updated at counteval == 0
42 chiN=N^0.5*(1-1/(4*N)+1/(21*N^2)); % expectation of
43 % ||N(0,I)|| == norm(randn(N,1))
44
45 % ----- Generation Loop -----
46
47 counteval = 0; % the next 40 lines contain the 20 lines of interesting code
48 while counteval < stopeval
49
50 % Generate and evaluate lambda offspring
51 for k=1:lambda,
52     arz(:,k) = randn(N,1); % standard normally distributed vector
53     arx(:,k) = xmean + sigma * (B*D * arz(:,k)); % add mutation % Eq. 37
54     arfitness(k) = feval(strfitnessfct, arx(:,k)); % objective function call
55     counteval = counteval+1;
56 end
57
58 % Sort by fitness and compute weighted mean into xmean
59 [arfitness, arindex] = sort(arfitness); % minimization
60 xmean = arx(:,arindex(1:mu))*weights; % recombination % Eq. 39
61 zmean = arz(:,arindex(1:mu))*weights; % == D^-1*B'*(xmean-xold)/sigma
62
63 % Cumulation: Update evolution paths
64 ps = (1-cs)*ps + (sqrt(cs*(2-cs)*mueff)) * (B * zmean); % Eq. 40
65 hsig = norm(ps)/sqrt(1-(1-cs)^(2*counteval/lambda))/chiN < 1.4+2/(N+1);
66 pc = (1-cc)*pc + hsig * sqrt(cc*(2-cc)*mueff) * (B*D*zmean); % Eq. 42

```

```

67
68 % Adapt covariance matrix C
69 C = (1-ccov) * C ... % regard old matrix % Eq. 43
70 + ccov * (1/mucov) * (pc*pc' ... % plus rank one update
71 + (1-hsig) * cc*(2-cc) * C) ...
72 + ccov * (1-l/mucov) ... % plus rank mu update
73 * (B*D*arz(:,arindex(1:mu))) ...
74 * diag(weights) * (B*D*arz(:,arindex(1:mu)))';
75
76 % Adapt step-size sigma
77 sigma = sigma * exp((cs/damps)*(norm(ps)/chiN - 1)); % Eq. 41
78
79 % Update B and D from C
80 if counteval - eigeneval > lambda/ccov/N/10 % to achieve O(N^2)
81 eigeneval = counteval;
82 C=triu(C)+triu(C,1)'; % enforce symmetry
83 [B,D] = eig(C); % eigen decomposition, B=normalized eigenvectors
84 D = diag(sqrt(diag(D))); % D contains standard deviations now
85 end
86
87 % Break, if fitness is good enough
88 if arfitness(1) <= stopfitness
89 break;
90 end
91
92 % Escape flat fitness
93 if arfitness(1) == arfitness(ceil(0.7*lambda))
94 sigma = sigma * exp(0.2+cs/damps);
95 disp('warning: flat fitness, consider reformulating the objective');
96 end
97
98 disp([num2str(counteval) ' : ' num2str(arfitness(1))]);
99
100 end % while, end generation loop
101
102 % ----- Ending Message -----
103
104 disp([num2str(counteval) ' : ' num2str(arfitness(1))]);
105 xmin = arx(:, arindex(1)); % Return best point of last generation.
106 % Notice that xmean is expected to be even
107 % better.
108
109 % -----
110 function f=felli(x)
111 N = size(x,1); if N < 2 error('dimension must be greater one'); end
112 f=1e6.^((0:N-1)/(N-1)) * x.^2; % condition number 1e6

```

C Reformulation of Learning Parameter c_{cov}

For sake of clarity, we consider to replace the learning coefficients in (43) in future.

$$\frac{c_{\text{cov}}}{\mu_{\text{cov}}} \quad \text{with} \quad c_1 \quad (53)$$

$$c_{\text{cov}} \left(1 - \frac{1}{\mu_{\text{cov}}}\right) \quad \text{with} \quad c_\mu \quad \text{and} \quad (54)$$

$$1 - c_{\text{cov}} \quad \text{with} \quad 1 - c_1 - c_\mu \quad , \quad (55)$$

and choosing (in replacing (49))

$$c_1 = \frac{2}{(n + 1.3)^2 + \mu_{\text{cov}}} \quad (56)$$

$$c_\mu = \min \left(2 \frac{\mu_{\text{cov}} - 2 + \frac{1}{\mu_{\text{cov}}}}{(n + 2)^2 + \mu_{\text{cov}}}, 1 - c_1 \right) \quad (57)$$

$$, \quad (58)$$

The resulting coefficients are quite similar to the previous. In contrast to the previous formulation, c_1 becomes monotonic in μ_{eff}^{-1} and $c_1 + c_\mu$ becomes virtually monotonic in μ_{eff} .

Another alternative is

$$c_1 = \frac{1}{m + 2\sqrt{m} + \frac{\mu_{\text{eff}}}{n}} \quad (59)$$

$$c_\mu = \min \left(1 - c_1, \frac{\alpha_\mu^0 + \mu_{\text{eff}} - 2 + \frac{1}{\mu_{\text{eff}}}}{m + 4\sqrt{m} + \frac{\mu_{\text{eff}}}{2}} \right) \quad (60)$$

$$\alpha_\mu^0 = 0.3 \quad , \quad (61)$$

where $m = \frac{n^2+n}{2}$ is the degrees of freedom in the covariance matrix. For $\mu_{\text{eff}} = 1$, the coefficient c_μ is now chosen to be larger than zero, as $\alpha_\mu^0 > 0$. Figure 7 compares the new learning rates with the old ones.

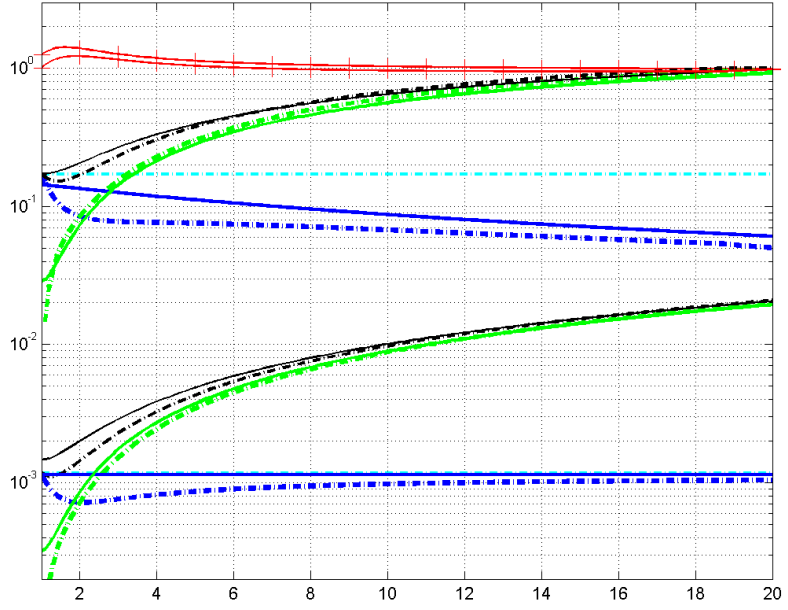
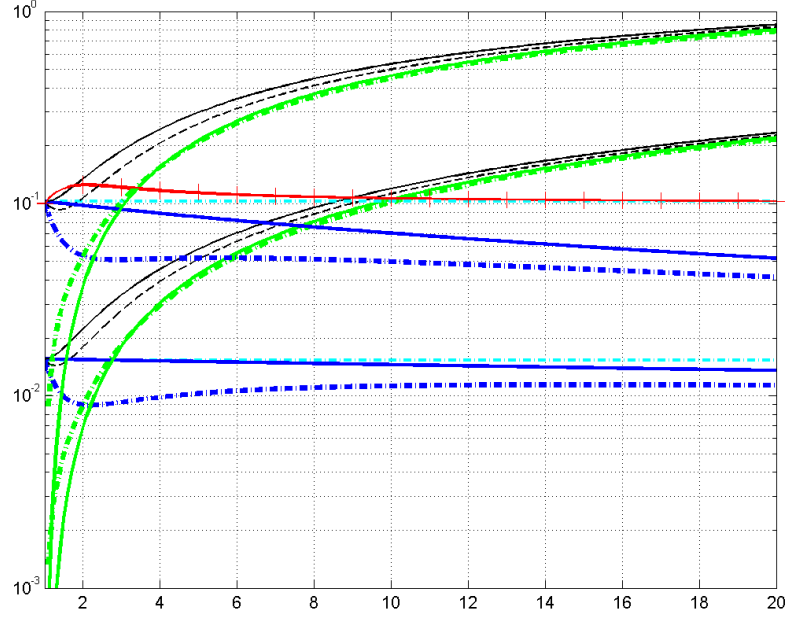


Figure 7: Learning rates c_1, c_μ (solid) and c_{cov} (dash-dotted) versus μ_{eff} . Above: Equations (56) etc. for $n = 3; 10$. Below: Equations (59) etc. for $n = 2; 40$. Black: $c_1 + c_\mu$ and c_{cov} ; blue: c_1 and $c_{\text{cov}}/\mu_{\text{cov}}$; green: c_μ and $(1 - 1/\mu_{\text{cov}})c_{\text{cov}}$; cyan: $2/(n^2 + \sqrt{2})$; red: $(c_1 + c_\mu)/c_{\text{cov}}$, above divided by ten. For $\mu_{\text{cov}} \approx 2$ the difference is maximal, because c_1 decreases much slower with increasing μ_{cov} and c_{cov} is non-monotonic in μ_{cov} (a main reason for the new formulation).