

# Adaptive Encoding for Optimization

Nikolaus Hansen

► **To cite this version:**

Nikolaus Hansen. Adaptive Encoding for Optimization. [Research Report] RR-6518, INRIA. 2008.  
<inria-00275983v3>

**HAL Id: inria-00275983**

**<https://hal.inria.fr/inria-00275983v3>**

Submitted on 29 Apr 2008

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

# *Adaptive Encoding for Optimization*

Nikolaus Hansen

**N° 6518**

April 2008

Thème COG



*Rapport  
de recherche*



## Adaptive Encoding for Optimization

Nikolaus Hansen\*

Thème COG — Systèmes cognitifs  
Équipes-Projets Adaptive Combinatorial Search et TAO

Rapport de recherche n° 6518 — April 2008 — 18 pages

**Abstract:** This report describes a general method for rendering search coordinate system independent, *Adaptive Encoding* (AE). Adaptive Encoding is applicable to any continuous domain search algorithm and includes (incremental) changes of the coordinate system, that is, changes of the representation of solutions. One attractive way to change the representation within AE is derived from the Covariance Matrix Adaptation (CMA). We prove that adaptive encoding recovers the CMA Evolution Strategy, when suitably applied to an evolution strategy with cumulative step-size adaptation. The proof implies that adaptive encoding provides the means to apply CMA-like representation changes to any search algorithm in continuous domain.

**Key-words:** optimization, evolutionary algorithms, encoding, representation, coordinate system, covariance matrix adaptation, invariance

\* Adaptive Combinatorial Search Team, Microsoft Research–INRIA Joint Centre. 28, rue Jean Rostand, 91893 Orsay Cedex, France. email:forename.name@inria.fr.

# Encodage adaptif pour optimisation

**Résumé :** Pas de résumé

**Mots-clés :** Pas de motclef

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Preliminary Notations and Definitions</b>	<b>4</b>
	Symbols . . . . .	5
<b>3</b>	<b>Adaptive Encoding</b>	<b>5</b>
<b>4</b>	<b>Example Applications of Adaptive Encoding</b>	<b>7</b>
	4.1 The Scaling Adaptive Evolutionary Algorithm . . . . .	7
	4.2 EDA (Estimation of Distribution Algorithm) . . . . .	9
<b>5</b>	<b>A Universal Update Rule: <math>\text{AE}_{\text{CMA}}</math></b>	<b>11</b>
	5.1 Choice of Parameters . . . . .	13
	5.2 $\text{AE}_{\text{CMA}}$ Recovers CMA-ES . . . . .	14
	5.3 Application of $\text{AE}_{\text{CMA}}$ -Update . . . . .	17
<b>6</b>	<b>Summary and Conclusions</b>	<b>17</b>

## 1 Introduction

In optimization or search, the problem encoding, that is the choice of the representation of the optimization problem is of utmost importance. A good representation, if available, can render any search problem trivial—finding a proper representation means essentially solving the problem. In an iterative search procedure, in principle, a good problem representation can be iteratively approached, just as a good solution to the problem is approached in the iteration sequence. Indeed, for example, variable metric methods like quasi-Newton methods [3, 4], covariance matrix adaptation (CMA) [7], or estimation of distribution algorithms (EDAs) [9] implicitly conduct a representational change. In evolutionary algorithms, given an additive (mutation) operator, a linear change of representation is equivalent with an appropriate linear transformation of the additive mutation [5]. Linear transformations of additive mutation operators, parameterized in step-sizes or covariance matrices, are well studied in evolutionary algorithms [1, 9].

In this report, we sketch an explicit framework for an iterative incremental representation change, adaptive encoding. The framework by itself is just about trivial. While the framework is quite general, this report considers only *linear* changes of the representation in continuous domain.

Searching for a linear representational change in continuous domain complements the original  $n$ -dimensional search problem with a second search problem of size  $n^2$ . The advantage from adaptive encoding is that these two search problems are decoupled. Consequently, any efficient representation change can be applied to any underlying search procedure.

In the next section we give the preliminary notations and definitions. In Section 3 the general idea of adaptive encoding is introduced. Section 4 explicates a few examples. In one of them adaptive encoding recovers an EDA. Section 5 proposes a universal, efficient update rule,  $\text{AE}_{\text{CMA}}$ , for the representation. We prove that  $\text{AE}_{\text{CMA}}$ -CSA-ES recovers CMA-ES.

## 2 Preliminary Notations and Definitions

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be an objective function to be minimized. Let a search algorithm propose new candidate solutions,  $\mathbf{x}$ , in an iterative procedure and typically evaluate them on  $f$ . Let further denote

$S$  the **state space** of the search algorithm;

$\mathcal{A} : S \rightarrow S$  an iteration step of the search algorithm  $\mathcal{A}$ ;

$T_B : S \rightarrow S$  an invertible transformation, the **encoding** of the state space or change of representation. The transformation  $T_B$  is here parameterized by a matrix  $\mathbf{B}$  and uniquely depends on the values in  $\mathbf{B}$ . The inverse  $T_B^{-1}$  denotes the decoding of the state space;

$\mathbf{B} \in \mathbb{R}^{n \times n}$  a full rank matrix. Columns of  $\mathbf{B}$  represent a new **coordinate system** and  $\mathbf{B}$  represents the respective a coordinate system transformation in  $\mathbb{R}^n$ . The linear *encoding* of candidate solutions reads  $\mathbf{B} : \mathbf{x} \mapsto \mathbf{B}\mathbf{x}$ ;

$\mathcal{U} : \mathbb{R}^{n \times n} \times S \rightarrow \mathbb{R}^{n \times n}$ ,  $(\mathbf{B}, s) \mapsto \mathcal{U}(\mathbf{B}, s) = \mathbf{B}^{\text{new}}$  the change of representation by updating the matrix  $\mathbf{B}$ . For convenience, we assume that all available and necessary information to update  $\mathbf{B}$  is included in algorithm state  $s$  and we will regularly write  $\mathcal{U}(\mathbf{B})$  instead of  $\mathcal{U}(\mathbf{B}, s)$ ;

From these definitions we first remark, that an iteration step of an algorithm can be surrounded by an encoding-decoding step according to

$$\mathcal{A}_B \equiv T_B \circ \mathcal{A} \circ T_B^{-1} , \quad (1)$$

defining a new algorithm  $\mathcal{A}_B : S \rightarrow S$ . If  $T_B$  is the identity the original algorithm  $\mathcal{A}$  is recovered.

By definition, *encoded* solutions (phenotypes) are represented in the *given* coordinate system, where also  $f$  is evaluated. Accordingly, the algorithm operates, by definition, on *decoded* solutions (genotypes). We assume, for convenience and w.l.o.g., that recently evaluated solutions are part of the algorithms state.

**Remark 1** (Evaluation of solutions). *In order to make use of Eq. (1), we have to ensure that candidate solutions are evaluated according to their original representation. Solutions must be encoded before evaluation. In other words,  $\mathcal{A}$  operates on  $f \circ \mathbf{B}$ .*

Considering Remark 1, we can execute the algorithm  $\mathcal{A}$  in any coordinate system of our choice. The new coordinate system, where the operations of  $\mathcal{A}$  are effectively conducted, is given by  $\mathbf{B}$ . Optimizing  $f \circ \mathbf{B}$  instead of  $f$  already renders  $\mathcal{A}$  independent of the given coordinate system (if  $\mathbf{B}$  is chosen independent of the given coordinate system). Eq. (1) becomes meaningful when we additionally adapt  $\mathbf{B}$ . Later we choose  $T_B$  such that changes of  $\mathbf{B}$  do not change the retained solutions in their original representation (encoded solutions).

Finally, we assume to have a performance measure when running an algorithm on a function.

**Remark 2** (Performance measure). *In the following, we assume to have a performance measure for running a search algorithm on the objective function  $f$ .*

The performance measure determines whether one algorithm is better than another. For example, a typical, quantitatively useful measure is the number of candidate solutions evaluated on  $f$  until a target function values is reached.

Further symbols which will be used repeatedly in this report are given in the next paragraph.

### Symbols

$\circ$  the composition of mappings, e.g.,  $(f \circ \mathbf{B})(\mathbf{x}) = f(\mathbf{B}(\mathbf{x}))$  for all  $\mathbf{x}$ .

$\leftarrow$  the assignment operation. We will write, for example,

$$\mathbf{B}\mathbf{B}^T \leftarrow \mathbf{C} ,$$

meaning we assign  $\mathbf{B}$  to a value such that  $\mathbf{B}\mathbf{B}^T = \mathbf{C}$ . If the assignment is not unique under the definition for  $\mathbf{B}$ , any compliant value is assigned.

$\mu, \lambda$  integers denoting the number of parental solutions and offspring solutions respectively.

$\mu_{\mathbf{w}} = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \geq 1$  the effective parent number, also denoted as *variance effective selection mass* [6].

$\mathbf{I}$  the identity  $n \times n$  unit matrix.

$i(f)$  the index of the  $i$ -the best solution in respect to function value  $f$ .

$\mathcal{N}_i(\mathbf{0}, \mathbf{C})$  a multi-variate normal distribution with zero mean and covariance matrix  $\mathbf{C}$ . The index  $i$  denotes a certain realization.

$R_i \in \mathbb{R}^n$  a random vector.

$\mathbf{B}^T$  the transpose of  $\mathbf{B}$ .

$w_i = \frac{\ln(\mu+1) - \ln i}{\mu \ln(\mu+1) - \sum_{j=1}^{\mu} \ln j}$ , for  $i = 1, \dots, \mu$ , positive weights summing to one.

## 3 Adaptive Encoding

Equation (1) represents an iteration step of a search algorithm with an additional encoding-decoding procedure. The encoding is, throughout this report, parameterized by a  $n \times n$ -matrix; it therefore introduces  $n^2$  additional degrees of freedom to the iteration step of the algorithm. Obviously, the idea is to find a “good” encoding for algorithm  $\mathcal{A}$ .

**Aim 1** (static encoding). *The goal of finding a good encoding is to find a transformation  $T_B$ , such that*

$$T_B \circ \mathcal{A} \circ T_B^{-1} \text{ outperforms } \mathcal{A} \text{ on } f$$



The static encoding is usually part of the design of the objective function. Equivalently, the algorithm can be modified specifically in regard to the given objective function (the encoding-decoding can certainly be interpreted as modification of the algorithm). The formalism of Aim 1 is not very interesting. To get a more interesting situation, we need to consider an *update* or *adaptation* of the encoding  $T_B$ .

**Definition 1.** (*Adaptive Encoding*) Given an algorithm,  $\mathcal{A}$ , an encoding,  $T_B$ , determined by  $\mathbf{B}$ , and an update,  $\mathcal{U}$ , the iteration step of an adaptively encoded algorithm in state  $s \in S$  is defined as

$$s \leftarrow T_B \circ \mathcal{A} \circ T_B^{-1}(s) \quad (2)$$

$$\mathbf{B} \leftarrow \mathcal{U}(\mathbf{B}, s) \quad (3)$$

where  $\leftarrow$  denotes the assignment operator and  $T_B \circ \mathcal{A} \circ T_B^{-1}(s) = T_B(\mathcal{A}(T_B^{-1}(s)))$ . We write

$$T_B \circ \mathcal{A} \circ T_B^{-1}; \mathcal{U}(T_B) \quad (4)$$

to denote the iteration step defined by Equations (2) and (3).

Obviously, any iterative algorithm  $\mathcal{A}$  can be plugged into the adaptive encoding mechanism.

**Proposition 1.** (*Adaptive Encoding is universal*) The Adaptive Encoding from Definition 1 can be applied to any search algorithm.

*Proof.* The proposition follows directly from the definition of  $T_B$  as invertible mapping from  $S$  to  $S$ .  $\square$

Even though Proposition 1 is just about trivial, it is of utmost importance for the implications of our results, because it establishes the general applicability of an effective adaptive encoding update rule.

Analogous to Aim 1, we consider the merits of an adaptive encoding.

**Aim 2** (adaptive encoding). Find an update  $\mathcal{U}$ , such that for a given  $T_0$  and a given (initial)  $T_B$

$$T_B \circ \mathcal{A} \circ T_B^{-1}; \mathcal{U}(T_B) \text{ outperforms } T_0 \circ \mathcal{A} \circ T_0^{-1} \text{ on } f.$$

The left iteration step updates the encoding, the right iteration step applies a constant encoding,  $T_0$ , to algorithm  $\mathcal{A}$ .

Taking only a single iteration step, Aim 2 does not depend on the update  $\mathcal{U}$  and it reduces to Aim 1. Consequently, Aim 2 becomes only interesting, when an iteration *sequence* is considered. Indeed, in a realistic automated scenario, Aim 2 can only be achieved in the iteration sequence.

Finally, we define three cases/scenarios when considering Aim 2.

**Scenario 1.** (*Standard scenario*) The initial  $T_B$  equals to  $T_0$ . Aim 2 shall be satisfied for most given  $T_0$ .

**Scenario 2.** (*Ambitious scenario*) The initial  $T_B$  equals to  $T_0$ . Aim 2 shall be satisfied for all given  $T_0$ .

Satisfying the ambitious scenario implies that no fixed optimal encoding  $T_B$  exists and a changing encoding can, in principle, be better than any fixed encoding. Both, the standard and the ambitious scenario are reasonable objectives, depending on the given objective function.

**Scenario 3.** (*Unrealistic scenario*) For all  $T_B \neq T_0$ , Aim 2 shall be satisfied.

As we might be able to choose  $T_B$  arbitrarily bad, it seems unrealistic that Aim 2 can be satisfied for any  $T_B \neq T_0$  in general.

The remainder of this report further investigates the idea of adaptive encoding as given in Definition 1.

## 4 Example Applications of Adaptive Encoding

In this section, we give two examples of how an evolutionary algorithm is reformulated as a simple base algorithm with adaptive encoding.

In order to define an adaptive encoding, we need to specify the encoding of the algorithms state space,  $T_B : S \rightarrow S$ , and the update of the encoding,  $\mathcal{U}$ . Both mappings depend, to a certain extend, on the baseline algorithm they are applied to.

### 4.1 The Scaling Adaptive Evolutionary Algorithm

We define in Algorithm 1 a “Baseline  $(\mu, \lambda)$ -EA” with  $\mu$  parental solutions, a recombination operator  $\mathbf{reco}$ ,  $\lambda$  offspring which are generated by adding a random vector  $R_i \in \mathbb{R}^n$ , and a selection operator. The  $(\mu, \lambda)$ -selection was chosen for notational simplicity. The results generalize to any selection operator that uses only the objective function values to take decisions (and not the parameter values of the candidate solutions itself).

The adaptive encoding of the baseline  $(\mu, \lambda)$ -EA, AE- $(\mu, \lambda)$ -EA, is given in Algorithm 2. It adds a decoding step in the beginning, encodes solutions for evaluation, adds an encoding step at the end and finally updates the encoding.

**Remark 3.** The encoding  $T_B$  in the adaptively encoded  $(\mu, \lambda)$ -EA (Algorithm 2) is of the form

$$T_B : (\mathbf{x}_1, \dots, \mathbf{x}_\mu, \sigma, y) \mapsto (\mathbf{B}\mathbf{x}_1, \dots, \mathbf{B}\mathbf{x}_\mu, \sigma, y) , \quad (5)$$

where  $y$  denotes remaining state variables.

---

#### Algorithm 1: Baseline $(\mu, \lambda)$ -EA

The  $R_i \in \mathbb{R}^n$  denote random vectors

---

- 1 let  $\mu, \lambda$  positive integers and  $\mu < \lambda$
  - 2 initialize  $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\} \in \mathbb{R}^n$  (a set of candidate solutions)
  - 3 **repeat**
  - 4    $\mathbf{x}_i = \mathbf{reco}(\mathbf{x}_1, \dots, \mathbf{x}_\mu) + R_i$ , for  $i = 1, \dots, \lambda$
  - 5    $f_i = f(\mathbf{x}_i)$ , for  $i = 1, \dots, \lambda$
  - 6    $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\} \leftarrow \mathbf{select}(\{\mathbf{x}_i, f_i\}_{i=1, \dots, \lambda})$
  - 7 **until** stopping criterion is met
-

We aim at investigating the relation of the adaptive encoding of the baseline  $(\mu, \lambda)$ -EA to conventional adaptations in evolutionary algorithms. Algorithm 3 complements the baseline  $(\mu, \lambda)$ -EA with an adaptive scaling of mutations. We will prove in the following that, given the appropriate  $T_B$  and  $\mathcal{U}$ , the adaptively encoded algorithm (Algorithm 2) is equivalent to the scaling adaptive Algorithm 3, where mutations are scaled according to a diagonal matrix  $\sigma$ .

**Assumption 1.** Let *select* pick solutions only depending on the function values (not depending on the actual solutions), and let *reco* commute with  $\mathbf{B}$ , in that

$$\mathbf{B} \cdot \text{reco}(\mathbf{x}_1, \dots, \mathbf{x}_\mu) = \text{reco}(\mathbf{B}\mathbf{x}_1, \dots, \mathbf{B}\mathbf{x}_\mu) \quad (6)$$

for all  $\mathbf{B}$  and all  $\mathbf{x}_1, \dots, \mathbf{x}_\mu$ .

**Theorem 1.** Given Assumption 1, the adaptively encoded  $(\mu, \lambda)$ -EA (Algorithm 2) implements the scaling adaptive  $(\mu, \lambda)$ -EA (Algorithm 3), if

$$\mathcal{U}(\sigma) = \text{adapt}_\sigma(\sigma) \quad (7)$$

---

**Algorithm 2:** Adaptively Encoded  $(\mu, \lambda)$ -EA

Shaded areas implement the adaptive encoding. The begin-end block implements the baseline  $(\mu, \lambda)$ -EA minimizing  $f \circ \mathbf{B}$ .

---

```

1 initialize  $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\} \in \mathbb{R}^n$  (a set of solutions)
2 initialize  $\mathbf{B} \in \mathbb{R}^{n \times n}$  (a transformation matrix)
3 repeat
4    $\mathbf{x}'_i = \mathbf{B}^{-1}\mathbf{x}_i$ , for  $i = 1, \dots, \mu$  // decode state
5   begin
6      $\mathbf{x}'_i = \text{reco}(\mathbf{x}'_1, \dots, \mathbf{x}'_\mu) + R_i$ , for  $i = 1, \dots, \lambda$ 
7      $f'_i = f \circ \mathbf{B}(\mathbf{x}'_i) = f(\mathbf{B}\mathbf{x}'_i)$ , for  $i = 1, \dots, \lambda$  // encode to eval
8      $\{\mathbf{x}'_1, \dots, \mathbf{x}'_\mu\} \leftarrow \text{select}(\{\mathbf{x}'_i, f'_i\}_{i=1, \dots, \lambda})$ 
9   end
10   $\mathbf{x}_i \leftarrow \mathbf{B}\mathbf{x}'_i$ , for  $i = 1, \dots, \mu$  // encode
11   $\mathbf{B} \leftarrow \mathcal{U}(\mathbf{B})$  // adapt encoding
12 until stopping criterion is met
```

---

**Algorithm 3:** Scaling Adaptive  $(\mu, \lambda)$ -EA

The shaded areas add to the baseline  $(\mu, \lambda)$ -EA

---

```

1 initialize  $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\} \in \mathbb{R}^n$  (a set of solutions)
2 initialize diagonal matrix  $\sigma$  (a scaling)
3 repeat
4    $\mathbf{x}_i = \text{reco}(\mathbf{x}_1, \dots, \mathbf{x}_\mu) + \sigma \times R_i$ , for  $i = 1, \dots, \lambda$ 
5    $f_i = f(\mathbf{x}_i)$ , for  $i = 1, \dots, \lambda$ 
6    $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\} \leftarrow \text{select}(\{\mathbf{x}_i, f_i\}_{i=1, \dots, \lambda})$ 
7    $\sigma \leftarrow \text{adapt}_\sigma(\sigma)$ 
8 until stopping criterion is met
```

---

holds with probability one.

*Proof.* Let initialize  $\mathbf{B} = \sigma$  in Algorithm 2 and let the initial  $\mathbf{x}_i$  be identical in both algorithms. We consider one iteration step of both algorithms.

First, we show that the same solutions are evaluated in both algorithms and consequently  $f'_i = f_i$ . That means, we have to prove that  $\mathbf{B}\mathbf{x}'_i$  in line 7 of Algorithm 2 equals to  $\mathbf{x}_i$  in line 5 of Algorithm 3, for all  $i = 1, \dots, \lambda$ .

$$\begin{aligned}
 \mathbf{B}\mathbf{x}'_i &= \mathbf{B}(\text{reco}(\mathbf{x}'_1, \dots, \mathbf{x}'_\mu) + R_i) \\
 &= \mathbf{B}\text{reco}(\mathbf{B}^{-1}\mathbf{x}_1, \dots, \mathbf{B}^{-1}\mathbf{x}_\mu) + \mathbf{B}R_i \\
 &= \text{reco}(\mathbf{B}\mathbf{B}^{-1}\mathbf{x}_1, \dots, \mathbf{B}\mathbf{B}^{-1}\mathbf{x}_\mu) + \sigma \times R_i \\
 &= \text{reco}(\mathbf{x}_1, \dots, \mathbf{x}_\mu) + \sigma \times R_i \\
 &= \mathbf{x}_i \quad \text{for } i = 1, \dots, \lambda
 \end{aligned} \tag{8}$$

Second, we prove that all  $\mathbf{x}_i$  are identical after one iteration step in both algorithms. According to Eq. (8),  $\mathbf{B}\mathbf{x}'_i$  in line 7 of Algorithm 2 equals to  $\mathbf{x}_i$  and therefore  $f'_i = f_i$ . Because `select` only depends on the function values, the individuals with the same indices are selected in both algorithms, respectively. Decoding in line 10 of Algorithm 2 assigns the desired equality.

Third, we remark that Eq. (7) implies that  $\mathbf{B} = \sigma$  remains valid after the iteration step.  $\square$

We have proven that under Assumption 1, adaptive encoding can equivalently replaced adaptively scaled mutations. We conclude by a few remarks on the applicability of this result.

**Remark 4.** *The update,  $\mathcal{U}$ , can always be chosen, such that Eq. (7) holds, if all necessary information is included in the algorithms state.*

**Remark 5.** *Some typical recombination operators [1] commute with linear transformations.*

**Remark 6.** *The diagonal matrix  $\sigma$  in the scaling adaptive algorithm resembles (overall) step-size adaptation, if all diagonal entries remain identical.*

## 4.2 EDA (Estimation of Distribution Algorithm)

Algorithm 4 implements a very simple Estimation of Distribution Algorithm (EDA) [9], denoted as sEDA in the following. Only the mean vector of the distribution is estimated for the next generation. Adaptive encoding for the sEDA is elementary and outlined in Algorithm 5. Only the mean vector  $\mathbf{m}$  needs to be transformed according to  $T_B^{-1}(\mathbf{m}) = \mathbf{B}^{-1}\mathbf{m}$ . The adaptation of the encoding is given in line 11 of Algorithm 5. The adaptation resembles the estimation of a covariance matrix in a common EDA. The common EDA is given in Algorithm 6, where the empirical covariance matrix is computed in line 7.

The next theorem states that the adaptively encoded sEDA unconditionally implements the common EDA.

**Theorem 2.** *The adaptively encoded sEDA (Algorithm 5) implements the common EDA (Algorithm 6).*

*Proof.* We assume that initially  $\mathbf{B}\mathbf{B}^\top = \mathbf{C}$  and that the initial  $\mathbf{m}$  is the same in both algorithms. We prove that the same holds true after one iteration step.

---

**Algorithm 4:** Simple EDA (sEDA)

$\mathcal{N}(\mathbf{0}, \mathbf{I})$  denotes a normal distribution with covariance matrix  $\mathbf{I}$ ;  $i(f)$  indicates the index of the  $i$ -th best solution

---

```

1 initialize  $\mathbf{m} \in \mathbb{R}^n$  (distribution mean)
2 repeat
3    $\mathbf{x}_i = \mathbf{m} + \mathcal{N}_i(\mathbf{0}, \mathbf{I})$ , for  $i = 1, \dots, \lambda$ 
4    $f_i = f(\mathbf{x}_i)$ , for  $i = 1, \dots, \lambda$ 
5    $\mathbf{m} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}_{i(f)}$ 
6 until stopping criterion is met
```

---



---

**Algorithm 5:** Adaptively Encoded sEDA

Shaded lines implement the adaptive encoding; the begin-end block embraces the sEDA;  $i(f)'$  indicates the index of the  $i$ -th best solution according to  $f'_i$

---

```

1 initialize  $\mathbf{m} \in \mathbb{R}^n$  (distribution mean)
2 initialize  $\mathbf{B} = \mathbf{I}$ 
3 repeat
4    $\mathbf{m}' = \mathbf{B}^{-1}\mathbf{m}$ 
5   begin
6      $\mathbf{x}'_i = \mathbf{m}' + \mathcal{N}(\mathbf{0}, \mathbf{I})$ , for  $i = 1, \dots, \lambda$ 
7      $f'_i = f(\mathbf{B}\mathbf{x}'_i)$ , for  $i = 1, \dots, \lambda$ 
8      $\mathbf{m}' \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}'_{i(f)'}$ 
9   end
10   $\mathbf{m} \leftarrow \mathbf{B}\mathbf{m}'$ 
11   $\mathbf{B}\mathbf{B}^\top \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{B}\mathbf{x}'_{i(f)'}) (\mathbf{B}\mathbf{x}'_{i(f)'})^\top$ 
12 until stopping criterion is met
```

---



---

**Algorithm 6:** Common EDA

Shaded areas add to the simple EDA

---

```

1 initialize  $\mathbf{m} \in \mathbb{R}^n$  (distribution mean)
2 initialize  $\mathbf{C} = \mathbf{I}$  (empirical covariance matrix)
3 repeat
4    $\mathbf{x}_i = \mathbf{m} + \mathcal{N}_i(\mathbf{0}, \mathbf{C})$ , for  $i = 1, \dots, \lambda$ 
5    $f_i = f(\mathbf{x}_i)$ , for  $i = 1, \dots, \lambda$ 
6    $\mathbf{m} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}_{i(f)}$ 
7    $\mathbf{C} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{x}_{i(f)} - \mathbf{m})(\mathbf{x}_{i(f)} - \mathbf{m})^\top$ 
8 until stopping criterion is met
```

---

First, we recognize that the same solutions are evaluated in both algorithms, because

$$\begin{aligned}
\mathbf{B}\mathbf{x}'_i &= \mathbf{B}(\mathbf{m}' + \mathcal{N}_i(\mathbf{0}, \mathbf{I})) && \text{in line 7 in Alg. 5} \\
&= \mathbf{B}\mathbf{B}^{-1}\mathbf{m} + \mathbf{B}\mathcal{N}_i(\mathbf{0}, \mathbf{I}) \\
&= \mathbf{m} + \mathcal{N}_i(\mathbf{0}, \mathbf{B}\mathbf{B}^T) \\
&= \mathbf{m} + \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \\
&= \mathbf{x}_i
\end{aligned} \tag{9}$$

The equation holds almost surely, if the random number realizations are chosen correspondingly in both algorithms. Consequently,  $f'_i = f_i$  and  $i(f) = i(f)'$  for all  $i = 1, \dots, \lambda$ . Therefore, the same solutions are used in the remaining operations and we find

$$\mathbf{m}_{\text{line 10}}^{\text{Alg. 5}} = \mathbf{B}\mathbf{m}' = \mathbf{B} \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}'_{i(f)'} = \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{B}\mathbf{x}'_{i(f)'} = \frac{1}{\mu} \sum_{i=1}^{\mu} \mathbf{x}_{i(f)} = \mathbf{m}_{\text{line 6}}^{\text{Alg. 6}}. \tag{10}$$

It remains to be shown that after the iteration step  $\mathbf{B}\mathbf{B}^T = \mathbf{C}$ . This follows directly from the choice of the update rule.

$$\begin{aligned}
\mathbf{B}\mathbf{B}^T &\leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{B}\mathbf{x}'_{i(f)'} - \mathbf{m})(\mathbf{B}\mathbf{x}'_{i(f)'} - \mathbf{m})^T \\
&= \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{x}_{i(f)} - \mathbf{m})(\mathbf{x}_{i(f)} - \mathbf{m})^T \\
&= \mathbf{C}
\end{aligned} \tag{11}$$

□

## 5 A Universal Update Rule: $\text{AE}_{\text{CMA}}$

In order to define an adaptive encoding, we needed to specify the encoding of the algorithms state space,  $T_B : S \rightarrow S$ , and the update rule,  $\mathcal{U}$ , of the encoding matrix  $\mathbf{B}$ . In this section, our aim is to derive a universal update  $\mathcal{U}$ , leaving only the choice of  $T_B$  as remaining, algorithm specific design issue.

Previously, we have used two update rules, both in order to recover a known evolutionary algorithm. In the first case, we had chosen  $\mathcal{U} = \text{adapt}_\sigma$  without further specification. In the second case, we had computed the decomposition of an empirical covariance matrix to derive  $\mathbf{B}$ . The latter depends in particular on a sufficiently large number of solutions  $\mu$ . In the following, we explicate a more universal update rule for the representation matrix  $\mathbf{B}$ . The update is derived from the equations for the covariance matrix update in the  $(\mu/\mu_w, \lambda)$ -CMA-ES [6, 7], denoted as  $\text{AE}_{\text{CMA}}$ , and explicated in Algorithm 7  $\text{AE}_{\text{CMA}}$ -Update.

All parameters and coefficients of  $\text{AE}_{\text{CMA}}$ -Update are in detail discussed in Section 5.1. The state variables are  $\mathbf{m}$ ,  $\mathbf{p}$  and  $\mathbf{C}$ . The mean  $\mathbf{m}$  is initialized to the mean of initial solutions of the search algorithm to which  $\text{AE}_{\text{CMA}}$ -Update is applied, and initially  $\mathbf{p} = \mathbf{0}$  and  $\mathbf{C} = \mathbf{I}$ . The covariance matrix update is based on differences between new solutions and the (former) mean.

**Proposition 2.** Let  $\sigma$  denote a step-size and  $\mu_w^{-1} = \sum_{i=1}^{\mu} w_i^2$ . Let  $\alpha_p = 1$ ,  $\alpha_0 = \frac{\sqrt{\mu_w}}{\sigma}$  and  $\alpha_i = \sigma^{-1}$ , for  $i = 1, \dots, \mu$ . Then, the algorithm  $AE_{CMA}$ -Update implements the update equations for the evolution path,  $\mathbf{p}$ , and the covariance matrix,  $\mathbf{C} = \mathbf{B}\mathbf{B}^T$ , in the  $(\mu/\mu_w, \lambda)$ -CMA-ES.

*Proof.* Assuming that  $\mathbf{x}_1, \dots, \mathbf{x}_\mu$  are the  $\mu$  best solutions in the recent iteration step, line 5 computes  $\mathbf{m}$  according to Eq.(3) in [6]. Lines 7 and 10 of Algorithm 7  $AE_{CMA}$ -Update replicate the covariance matrix update equations (17) and (22) in [6] with added or renamed normalization coefficients, denoted as  $\alpha$ . Substituting the coefficients as given results in the original equations.  $\square$

The Algorithm 7  $AE_{CMA}$ -Update implements the covariance matrix update of CMA-ES with additional coefficients  $\alpha$  to be specified (see Section 5.1). In the context of CMA-ES, this update works well for any choice of  $\mu$  [6].

Depending on the application of  $AE_{CMA}$ -Update, a slow change of  $\mathbf{B}$  might be desirable (see also Section 5.3). While  $\mathbf{C}$  will only change slowly, as long as  $c_1$  and  $c_\mu$  are small, the decomposition of  $\mathbf{C}$  does not ensure a similar behavior for  $\mathbf{B}^\circ$  and  $\mathbf{D}$ . For this reason, the diagonal elements are sorted in  $\mathbf{D}$ . We conjecture, that lines 11 to 13 can be replaced by a cholesky decomposition, or, more promising, by an incremental cholesky update, similar to [8], that ensures small changes, as long as  $c_1$  and  $c_\mu$  are small. In this case, it might be sufficient to only encode the solutions for the function evaluation and, as an approximation, completely abandon the encoding-decoding of the algorithms state.

In the next section, we will discuss the choice of the coefficients and of the remaining parameters of the procedure. Section 5.2 shows how  $AE_{CMA}$ -Update recovers CMA-ES.

---

**Algorithm 7:**  $AE_{CMA}$ -Update( $\{\mathbf{x}_1, \dots, \mathbf{x}_\mu\}$ )

updates the encoding matrix  $\mathbf{B}$  using the  $\mu$  recent best-ranked candidate solutions

---

- 1 given parameters  $w_i, c_p, c_1, c_\mu$  from Section 5.1
  - 2 given  $\mathbf{m} \in \mathbb{R}^n$ ,  $\mathbf{p} \in \mathbb{R}^n$  and  $\mathbf{C} \in \mathbb{R}^{n \times n}$  from last iteration
  - 3 let matrices  $\mathbf{B}^\circ$  orthogonal, and  $\mathbf{D}$  diagonal, with diagonal elements sorted in ascending order
  - 4  $\mathbf{m}^- = \mathbf{m}$
  - 5  $\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_i$  // Eq.(3) in [6]
  - 6 set scalars  $\alpha_i \geq 0$ , for  $i = 0, \dots, \mu$ , cf. Sect. 5.1
  - 7  $\mathbf{p} \leftarrow (1 - c_p) \mathbf{p} + \sqrt{c_p(2 - c_p)} \alpha_0 (\mathbf{m} - \mathbf{m}^-)$  // Eq.(17) in [6]
  - 8  $\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \alpha_i^2 (\mathbf{x}_i - \mathbf{m}^-)(\mathbf{x}_i - \mathbf{m}^-)^T$  // rank- $\mu$  matrix
  - 9 set scalar  $\alpha_p \geq 0$ , cf. Sect. 5.1
  - 10  $\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \alpha_p \mathbf{p}\mathbf{p}^T + c_\mu \mathbf{C}_\mu$  // Eq.(22) in [6]
  - 11  $\mathbf{B}^\circ \mathbf{D} \mathbf{D} \mathbf{B}^\circ \leftarrow \mathbf{C}$  // eigendecomposition
  - 12 optionally normalize  $\mathbf{D}$
  - 13  $\mathbf{B} \leftarrow \mathbf{B}^\circ \mathbf{D}$  // encoding matrix
-

## 5.1 Choice of Parameters

In Algorithm 7, the scalars  $\alpha_p$  and  $\alpha_i$  for  $i = 0, \dots, \mu$ , need to be chosen. They normalize the input entries for the covariance matrix update (most of them are the difference between a new solution and the former mean). In the original CMA, we can derive the expected lengths of the input entries from the sampling procedure. Under random selection the normalized entries are distributed according to

$$\mathcal{N}(\mathbf{0}, \mathbf{C}) = \mathbf{B} \mathcal{N}(\mathbf{0}, \mathbf{I}) . \quad (12)$$

In general, we cannot assume to know the expected lengths of the input entries, therefore we need to normalize them. In Eq. (12), the expected squared length of the decoded input entry,  $E\|\mathbf{B}^{-1}\mathbf{B}\mathcal{N}(\mathbf{0}, \mathbf{I})\|^2$ , computes to  $n$  suggesting a normalization to length  $\sqrt{n}$ . Keeping this in mind, we discuss the choice of the scalar coefficients in turn.<sup>1</sup>

$\alpha_0 = \frac{\sqrt{n}}{\|\mathbf{B}^{-1}(\mathbf{m} - \mathbf{m}^-)\|}$ , normalizes the difference  $\mathbf{B}^{-1}(\mathbf{m} - \mathbf{m}^-)$  to length  $\sqrt{n}$ . Consequently, only the direction is relevant and the absolute size of the difference is disregarded.

$\alpha_i = \frac{\sqrt{n}}{\|\mathbf{B}^{-1}(\mathbf{x}_i - \mathbf{m}^-)\|}$ , for  $i = 1, \dots, \mu$ , is the conservative choice, where the length of the difference  $\mathbf{B}^{-1}(\mathbf{x}_i - \mathbf{m}^-)$  is disregarded. In general, we recommend to choose

$$\alpha_i = \frac{\sqrt{n}}{\max\left(\frac{l_i}{\beta}, \text{median}(l_j)_{j=1, \dots, \mu}\right)} \quad \text{for } i = 1, \dots, \mu , \quad (13)$$

where  $l_i = \|\mathbf{B}^{-1}(\mathbf{x}_i - \mathbf{m}^-)\|$ . In this way, the median is set to “length”  $\sqrt{n}$  and the maximal length is set to  $\beta\sqrt{n}$  with  $\beta \geq 1$ . We recommend  $\beta = 2$ . Unusual large entries may, for example, occur if solutions are originally sampled from a distribution with heavy tails. By chance, an outranging solution could enter the procedure despite a bad objective function value and an unjustified very large change of  $\mathbf{B}$  would result.

$\alpha_p = 1$  will be the usual choice, while  $\alpha_p = \frac{\sqrt{n}}{\|\mathbf{p}\|}$  is a conservative alternative and will not allow to utilize the evolution path effectively;  $\alpha_p = 0$  would be even more conservative.

Finally, we give the default settings for the constants used in  $\text{AE}_{\text{CMA-Update}}$  and discuss the choices in turn.

$c_p = \frac{1}{\sqrt{n}}$  is the learning constant for the evolution path, which should be usually between  $\frac{1}{\sqrt{n}}$  and  $\frac{2}{n+1}$  [7]. For larger  $c_p$ , the effect of the evolution path will attenuate. The backward time horizon for the evolution path is roughly  $c_p^{-1}$ . We choose as default the “conservative” limit of the useful range, *i.e.* a comparatively large  $c_p$ . The most conservative choice would be  $c_p = 1$ .

<sup>1</sup>We ignore the case of denominators being zero, where the respective coefficient  $\alpha$  can be set to any positive number.



$w_i = \frac{\ln(\mu+1) - \ln i}{\mu \ln(\mu+1) - \sum_{j=1}^{\mu} \ln j}$ , for  $i = 1, \dots, \mu$  are the recombination weights. They (must) sum to one and obey  $w_1 \geq \dots \geq w_{\mu} \geq 0$ . Generally, we choose  $\mu$  being half of the overall generated number of solutions per iteration (before selection).

$c_1 = \alpha_c \frac{0.2}{(n+1.3)^2 + \mu_w}$  is the learning rate for the rank-one update in line 10 of  $\text{AE}_{\text{CMA}}\text{-Update}$  (middle summand), with  $\alpha_c = 1$  as default. The denominator being quadratic in  $n$  reflects the degrees of freedom in the encoding matrix  $\mathbf{B}$ . The formula is derived as a simplification from the original formulation in [6].

$c_{\mu} = \alpha_c 0.2 \frac{\mu_w - 2 + \frac{1}{\mu_w}}{(n+2)^2 + \alpha_{\mu} \mu_w}$  is the learning rate for the rank- $\mu$  update (right summand in line 10), with  $\alpha_c = 1$  and  $\alpha_{\mu} = 0.2$  as default. With increasing  $\mu_w$ , the learning rate increases and gets close to one.

$\alpha_c = 1$  must be chosen positive and such that  $c_1 + c_{\mu} \leq 1$ . The default value of one is about ten times smaller, *i.e.* considerably more conservative, than for CMA-ES. Too large values for  $\alpha_c$  potentially lead to a failure. Too small values slow down the adaptation. In any case, at least a minimalistic parameter study for  $\alpha_c$  is recommended.

The final parameter setting needs to be decided specifically for a given algorithm. We believe that the given guidelines will be usually sufficient to find good settings with reasonable effort. For  $\text{AE}_{\text{CMA}}\text{-CSA-ES}$  (meaning CMA-ES), a good setting works across many objective functions and the identification needed to be conducted only once on a few simple test functions.

## 5.2 $\text{AE}_{\text{CMA}}$ Recovers CMA-ES

The CMA-ES implements two separate adaptation mechanisms, one for step-size control and another for covariance matrix adaptation. Our next application of  $\text{AE}_{\text{CMA}}$  is based on cumulative step-size adaptation (CSA, also denoted as *path length control*). The  $(\mu/\mu_w, \lambda)$ -CSA-ES is explicated in Algorithm 8. Its default parameters are

$$c_{\sigma} = \frac{\mu_w + 2}{n + \mu_w + 3} \quad \text{and} \quad d_{\sigma} = 1 + 2 \max \left( 0, \sqrt{\frac{\mu_w - 1}{n + 1}} - 1 \right) + c_{\sigma} . \quad (14)$$

For applying adaptive encoding, we choose the following invertible encoding for the state variables in CSA-ES,

$$T_B : (\mathbf{m}, \mathbf{p}_{\sigma}, \sigma) \mapsto (\mathbf{B}\mathbf{m}, \mathbf{B}^{\circ}\mathbf{p}_{\sigma}, \sigma) . \quad (15)$$

The  $\text{AE}_{\text{CMA}}\text{-}(\mu/\mu_w, \lambda)\text{-CSA-ES}$  applies adaptive encoding according to Algorithm 7,  $\text{AE}_{\text{CMA}}\text{-Update}$ , to the  $(\mu/\mu_w, \lambda)$ -CSA-ES and is given in Algorithm 9. The  $\mu$  best (encoded) solutions are used as input to  $\text{AE}_{\text{CMA}}\text{-Update}$  (line 18 in Algorithm 9). The encoding  $T_B$  solely depends on the matrix  $\mathbf{B}$ , as  $\mathbf{B}^{\circ}$  can be computed from  $\mathbf{B}$  by normalizing its columns to length one.

**Theorem 3** (Recovery of CMA-ES). *Given  $T_B$  as in Eq. (15) and the scalars for Algorithm 7  $\text{AE}_{\text{CMA}}\text{-Update}$  in each iteration as given in Proposition 2, then the  $\text{AE}_{\text{CMA}}\text{-}(\mu/\mu_w, \lambda)\text{-CSA-ES}$  (Algorithm 9) implements the  $(\mu/\mu_w, \lambda)\text{-CMA-ES}$ .*

*Proof.* As in our previous proofs, we assume the same initial state for  $\text{AE}_{\text{CMA}}(\mu/\mu_w, \lambda)$ -CSA-ES and  $(\mu/\mu_w, \lambda)$ -CMA-ES and first consider the sampled and

---

**Algorithm 8: CSA-ES**


---

```

1 initialize  $\mathbf{m} \in \mathbb{R}^n$  (distribution mean)
2 initialize  $\mathbf{p}_\sigma = \mathbf{0}$  (evolution path)
3 initialize  $\sigma > 0$  (step-size)
4 repeat
5    $\mathbf{x}_i = \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{I})$ , for  $i = 1, \dots, \lambda$ 
6    $f_i = f(\mathbf{x}_i)$ , for  $i = 1, \dots, \lambda$ 
7    $\mathbf{m}^- = \mathbf{m}$ 
8    $\mathbf{m} \leftarrow \sum_{i=1}^\mu w_i \mathbf{x}_{i(f)}$ 
9    $\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{c_\sigma (2 - c_\sigma) \mu_w} \frac{1}{\sigma} (\mathbf{m} - \mathbf{m}^-)$ 
10   $\sigma \leftarrow \sigma \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right)$ 
11 until stopping criterion is met

```

---



---

**Algorithm 9:  $\text{AE}_{\text{CMA}}$ -CSA-ES**


---

Shaded lines implement the adaptive encoding,  $\text{AE}_{\text{CMA}}$ , including the update of  $\mathbf{B}$  and  $\mathbf{B}^\circ$ . The begin-end block marks the original CSA-ES

---

```

1 initialize  $\mathbf{m} \in \mathbb{R}^n$  (distribution mean)
2 initialize  $\mathbf{p}_\sigma = \mathbf{0}$  (evolution path)
3 initialize  $\sigma > 0$  (step-size)
4 initialize  $\mathbf{B} = \mathbf{B}^\circ = \mathbf{I}$ 
5 repeat
6    $\mathbf{m}' = \mathbf{B}^{-1} \mathbf{m}$ 
7    $\mathbf{p}'_\sigma = \mathbf{B}^{\circ\top} \mathbf{p}_\sigma$ 
8   begin
9      $\mathbf{x}'_i = \mathbf{m}' + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{I})$ , for  $i = 1, \dots, \lambda$ 
10     $f'_i = f(\mathbf{B} \mathbf{x}'_i)$ , for  $i = 1, \dots, \lambda$ 
11     $\mathbf{m}'^- = \mathbf{m}'$ 
12     $\mathbf{m}' \leftarrow \sum_{i=1}^\mu w_i \mathbf{x}'_{i(f)'}$ 
13     $\mathbf{p}'_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}'_\sigma + \sqrt{c_\sigma (2 - c_\sigma) \mu_w} \frac{1}{\sigma} (\mathbf{m}' - \mathbf{m}'^-)$ 
14     $\sigma \leftarrow \sigma \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}'_\sigma\|}{E\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right)$ 
15  end
16   $\mathbf{m} = \mathbf{B} \mathbf{m}'$ 
17   $\mathbf{p}_\sigma = \mathbf{B}^\circ \mathbf{p}'_\sigma$ 
18   $\text{AE}_{\text{CMA}}\text{-Update}(\{\mathbf{B} \mathbf{x}'_1, \dots, \mathbf{B} \mathbf{x}'_\mu\})$ 
19 until stopping criterion is met

```

---

evaluated solutions  $\mathbf{B}\mathbf{x}'_i$ , for  $i = 1, \dots, \lambda$ , in Algorithm 7. We have

$$\begin{aligned} \mathbf{B}\mathbf{x}'_i &= \mathbf{B}(\mathbf{m}' + \sigma\mathcal{N}_i(\mathbf{0}, \mathbf{I})) \\ &= \mathbf{B}\mathbf{m}' + \sigma\mathbf{B}\mathcal{N}_i(\mathbf{0}, \mathbf{I}) \\ &= \mathbf{m} + \sigma\mathcal{N}_i(\mathbf{0}, \mathbf{B}\mathbf{B}^\top) \\ &= \mathcal{N}_i(\mathbf{m}, \sigma^2\mathbf{C}) \quad , \end{aligned} \tag{16}$$

which is the equation for generating new solutions in CMA-ES (Eq. (2) in [6]). Consequently, the same (encoded) solutions are generated (given respective random number realizations) and evaluated in  $\text{AE}_{\text{CMA-CSA-ES}}$  and CMA-ES. Also  $i(f) = i(f)'$  and the corresponding solutions are selected for the remaining updates.

We need to show that  $\text{AE}_{\text{CMA-CSA-ES}}$  recovers the update of all five state variables of CMA-ES,  $\mathbf{m}$ ,  $\mathbf{p}_\sigma$ ,  $\sigma$ ,  $\mathbf{p}$  and  $\mathbf{C}$ . We treat each variable in turn.

The new mean in  $\text{AE}_{\text{CMA-CSA-ES}}$  obeys

$$\begin{aligned} \mathbf{m}_{\text{line 16}}^{\text{Alg. 9}} &= \mathbf{B}\mathbf{m}' = \mathbf{B} \sum_{i=1}^{\mu} w_i \mathbf{x}'_{i(f)'} = \sum_{i=1}^{\mu} w_i \mathbf{B}\mathbf{x}'_{i(f)'} \\ &= \sum_{i=1}^{\mu} w_i \mathbf{x}_{i(f)} \quad , \end{aligned} \tag{17}$$

which is the equation for updating the mean in CMA-ES (Eq. (3) in [6]). Next, we investigate the evolution path for the step-size, starting from line 17.

$$\begin{aligned} \mathbf{p}_\sigma &= \mathbf{B}^\circ \mathbf{p}'_\sigma \\ &\leftarrow \mathbf{B}^\circ \left( (1 - c_\sigma) \mathbf{p}'_\sigma + \sqrt{c_\sigma(2 - c_\sigma)\mu_w} \frac{1}{\sigma} (\mathbf{m}' - \mathbf{m}'^-) \right) \\ &= (1 - c_\sigma) \mathbf{B}^\circ \mathbf{p}'_\sigma + \sqrt{c_\sigma(2 - c_\sigma)\mu_w} \frac{1}{\sigma} \mathbf{B}^\circ (\mathbf{m}' - \mathbf{m}'^-) \\ &= (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)\mu_w} \frac{1}{\sigma} \mathbf{B}^\circ \mathbf{B}^{-1} (\mathbf{m} - \mathbf{m}^-) \end{aligned} \tag{18}$$

We compute the rightmost term to

$$\begin{aligned} \mathbf{B}^\circ \mathbf{B}^{-1} (\mathbf{m} - \mathbf{m}^-) &= \mathbf{B}^\circ (\mathbf{B}^\circ \mathbf{D})^{-1} (\mathbf{m} - \mathbf{m}^-) \\ &= \mathbf{B}^\circ \mathbf{D}^{-1} \mathbf{B}^{\circ\top} (\mathbf{m} - \mathbf{m}^-) \quad , \end{aligned} \tag{19}$$

where  $\mathbf{B} = \mathbf{B}^\circ \mathbf{D}$  holds according to Algorithm 7. Equations (18) and (19) recover the update rule for the evolution path  $\mathbf{p}_\sigma$  in CMA-ES (Eq. (23) in [6]).

The update of step-size  $\sigma$  computes  $\|\mathbf{p}'_\sigma\| = \|\mathbf{B}^\circ \mathbf{p}_\sigma\|$ . Because  $\mathbf{B}^\circ$  is orthogonal we have  $\|\mathbf{B}^\circ \mathbf{p}_\sigma\| = \|\mathbf{p}_\sigma\|$ . Consequently, the step-size update is identical in  $\text{AE}_{\text{CMA-CSA-ES}}$  and CMA-ES.

Finally, completing the proof,  $\mathbf{p}$  and  $\mathbf{C}$  are updated in Algorithm  $\text{AE}_{\text{CMA-Update}}$  according to the CMA-ES, see Proposition 2.  $\square$

Theorem 3 supports the hypothesis that  $\text{AE}_{\text{CMA-Update}}$  is an efficient way to update the representation matrix  $\mathbf{B}$ , as CMA-ES is known to efficiently adapt the principle axes of the coordinate system, where the independent sampling takes place. In the next section, another application of  $\text{AE}_{\text{CMA-Update}}$  is sketched.

### 5.3 Application of $\text{AE}_{\text{CMA}}$ -Update

The scaling adaptive  $(\mu, \lambda)$ -EA, Algorithm 3, samples new solutions without dependencies between variables in the given coordinate system, because  $\sigma$  is a diagonal matrix. Rendering the scaling adaptive  $(\mu, \lambda)$ -EA coordinate system independent means that the distribution can reveal correlations with respect to the given coordinate system (if  $\sigma \neq \mathbf{I}$  is not the identity).

For the scaling adaptive  $(\mu, \lambda)$ -EA the encoding

$$T_B : (\mathbf{x}_1, \dots, \mathbf{x}_\mu, \sigma, y) \mapsto (\mathbf{B}\mathbf{x}_1, \dots, \mathbf{B}\mathbf{x}_\mu, \sigma, y) \quad (20)$$

suggests itself, where  $y$  denotes all further state variables.

The step-size matrix  $\sigma$  is not transformed. An appropriate mapping for a covariance matrix  $\sigma^2 \mapsto \mathbf{B}\sigma^2\mathbf{B}^T$  would not preserve the diagonal property. A transformation of the diagonal of  $\sigma$  with  $\mathbf{B}$  can lead, by chance, to very small entries—a very undesirable result for a step-size. Because  $\sigma$  is not encoded, it is important that changes of  $\mathbf{B}$  are modest.

A vector  $\mathbf{r}$ , that denotes a direction rather than a solution point in state space (as for example a velocity in Particle Swarm Optimization [2]), would arguably map to  $\mathbf{B}\mathbf{r}$ .

Using  $\mathbf{B}^\circ$  instead of  $\mathbf{B}$  in Eq. (20) is a possible alternative. Then, the step-size matrix  $\sigma$  needs to learn the scaling which can be otherwise provided by the diagonal matrix  $\mathbf{D}$ .

We experimented using the Cauchy distribution for the random vector  $R_i$  in Algorithm 3 as a rather different example compared to CSA. The adaptive encoding was successfully applied in this case and made the Baseline  $(\mu, \lambda)$ -EA roughly thousand times faster on proto-typical non-separable problems.

## 6 Summary and Conclusions

We have outlined an *adaptive* change of representation in continuous domain search, denoted as *adaptive encoding* (AE). The idea is simple: after each iteration step, (i) the algorithms state is “encoded”, (ii) the encoding is adapted, and (iii) the algorithms state is “decoded” again for the next iteration, using the updated inverse encoding. Further, candidate solutions are encoded for their evaluation on the objective function, respectively. The implications from this simple procedure are surprisingly far-reaching.

1. Well-known evolutionary algorithms can be recovered from basic algorithms in that a specific adaptive encoding is applied.
2. The update of the covariance matrix in the CMA-ES can be entirely formulated by adaptive encoding, as adaptation of a representation matrix  $\mathbf{B}$  (Proposition 2). The respective update rule is given in  $\text{AE}_{\text{CMA}}$ -Update (Algorithm 7). Applied to the cumulative step-size adaptation (CSA) evolution strategy, the  $\text{AE}_{\text{CMA}}$  can recover the CMA-ES (Theorem 3).
3. Adaptive encoding can be applied to any search algorithm. Consequently, with the  $\text{AE}_{\text{CMA}}$ , the “covariance matrix update” of CMA becomes applicable to any continuous domain search algorithm. The  $\text{AE}_{\text{CMA}}$  can render any search algorithm independent of the coordinate system, in particular

rotationally invariant. We anticipate successful applications of  $\text{AE}_{\text{CMA}}$  primarily to population based, stochastic algorithms.

## References

- [1] H.-G. Beyer and H.-P. Schwefel. Evolution strategies: A comprehensive introduction. *Natural Computing*, 1(1):3–52, 2002.
- [2] M. Clerc and J. Kennedy. The particle swarm-explosion, stability, and convergence in a multidimensional complex space. *Evolutionary Computation, IEEE Transactions on*, 6(1):58–73, 2002.
- [3] W. Davidon. Variable Metric Method for Minimization. *SIAM Journal on Optimization*, 1:1, 1991.
- [4] D. Goldfarb. A Family of Variable-Metric Methods Derived by Variational Means. *Mathematics of Computation*, 24(109):23–26, 1970.
- [5] N. Hansen. Invariance, self-adaptation and correlated mutations in evolution strategies. In M. Schoenauer, K. Deb, G. Rudolph, X. Yao, E. Lut-ton, J. Merelo, and H.-P. Schwefel, editors, *Parallel Problem Solving from Nature—PPSN VI, Proceedings*, pages 355–364, Paris, 2000. Springer, Berlin.
- [6] N. Hansen. The CMA evolution strategy: a comparing review. In J. Lozano, P. Larrañaga, I. Inza, and E. Bengoetxea, editors, *Towards a new evolutionary computation. Advances on estimation of distribution algorithms*, pages 75–102. Springer, 2006.
- [7] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.
- [8] C. Igel, T. Suttorp, and N. Hansen. A computational efficient covariance matrix update and a (1+1)-CMA for evolution strategies. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2006)*, pages 453–460. ACM Press, 2006.
- [9] P. Larrañaga. A review on estimation of distribution algorithms. In P. Larrañaga and J. Lozano, editors, *Estimation of distribution algorithms*, pages 80–90. Kluwer Academic Publishers, 2002.



---

Centre de recherche INRIA Saclay – Île-de-France  
Parc Orsay Université - ZAC des Vignes  
4, rue Jacques Monod - 91893 Orsay Cedex (France)

Centre de recherche INRIA Bordeaux – Sud Ouest : Domaine Universitaire - 351, cours de la Libération - 33405 Talence Cedex  
Centre de recherche INRIA Grenoble – Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier  
Centre de recherche INRIA Lille – Nord Europe : Parc Scientifique de la Haute Borne - 40, avenue Halley - 59650 Villeneuve d'Ascq  
Centre de recherche INRIA Nancy – Grand Est : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex  
Centre de recherche INRIA Paris – Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex  
Centre de recherche INRIA Rennes – Bretagne Atlantique : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex  
Centre de recherche INRIA Sophia Antipolis – Méditerranée : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399