# Dynamic Problem Encoding for Optimization

Nikolaus Hansen

May 7, 2008

## Content

Special thanks to
Anne Auger
Raymond Ros
Marc Schoenauer
Olivier Teytaud

*Einstein once spoke of the "unreasonable effectiveness of mathematics" in describing how the natural world works. Whether one is talking about basic physics, about the increasingly important environmental sciences, or the transmission of disease, mathematics is never any more, or any less, than a way of thinking clearly. As such, it always has been and always will be a valuable tool, but only valuable when it is part of a larger arsenal embracing analytic experiments and, above all, wide-ranging imagination.*

Lord Kay

## Problem Statement: Search

Continuous Domain Search/Optimization

- Task: **minimize** a **objective function** (*fitness* function, *loss* function) in **continuous** domain

$$f : \mathcal{X} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \boldsymbol{x} \mapsto f(\boldsymbol{x})$$

- **Black Box** scenario (direct search scenario)



  - gradients are not available or not useful
  - problem domain specific knowledge is used only within the black box, e.g. within an appropriate encoding

- Search costs: number of function evaluations

# Problem Statement and Objectives

Continuous Domain Search/Optimization

- Goal
  - ▸ fast convergence toward the global optimum

    . . . or to a robust solution $x$
  - ▸ solution $x$ with **small function value** with **least search cost**

    there are two (conflicting) objectives

- Typical Examples
  - ▸ shape optimization (e.g. using CFD)         curve fitting, airfoils
  - ▸ parameter calibration         controller, plants, images
  - ▸ model calibration         biological, physical

- Difficulties
  - ▸ exhaustive search is infeasible
  - ▸ deterministic search is often not successful
  - ▸ (naive) random search takes too long

**Approach**: stochastic search, Evolutionary Algorithms

. . . Interface to real world problems

**Problem Statement**
○○●

The Challenges
○○○

Evolution Strategy
○○○○○○○○

CMA
○○○○○○

Evaluation
○○○○○

Adaptive Encoding
○○○○○○○

# Problem Formulation

A real world problem requires

- a representation; the encoding of problem parameters into $x \in \mathcal{X} \subset \mathbb{R}^n$
- the definition of a objective function $f : x \mapsto f(x)$ to be minimized

One might distinguish two approaches

## Natural Encoding

Use a "natural" encoding and **design the optimizer** with respect to the problem e.g. use of specific "genetic operators"

frequently done in discrete domain

## Concerned Encoding (Pure Black Box)

Use problem specific knowledge for encoding and use a **"generic" optimizer**

frequently done in continuous domain

Advantage: Sophisticated and well-validated optimizers can be used

How about *Adaptive Encoding*?

...function properties

## Objective Function Properties

We assume $f : \mathcal{X} \subset \mathbb{R}^n \to \mathbb{R}$ to have at least moderate dimensionality, say $n \not\ll 10$, and to be *non-linear, non-convex, and non-separable*.

Additionally, $f$ can be

- multimodal

  there are eventually many local optima

- non-smooth

  derivatives do not exist

- discontinuous
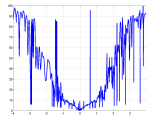- ill-conditioned
- noisy
- . . .

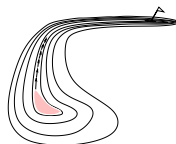     **Goal** : cope with any of these function properties

there are related to real-world problems

# What Makes a Function Difficult to Solve?

Why stochastic search?

- ruggedness
  non-smooth, discontinuous, multimodal, and/or
  noisy function

- dimensionality

  (considerably) larger than three

- non-separability
  dependencies between the objective variables

- ill-conditioning



cut from 5-D solvable example



a narrow ridge

# How Can a Difficult Function Be Solved?

. . . therefore. . .

| The Problem | What can be done |
|---|---|
| Ruggedness | **non-local** policy, large sampling width (step-size) <br> <span style="color:blue">as large as possible while preserving a reasonable convergence speed</span> |
| | stochastic, non-elitistic, **population-based** method <br> recombination operator <br> <span style="color:blue">serves as repair mechanism</span> |
| Dimensionality, Non-Separability | exploiting the problem structure <br> <span style="color:blue">locality, neighborhood, encoding</span> |
| Ill-conditioning | second order approach <br> <span style="color:blue">changes the neighborhood metric</span> |

## Stochastic Search

A black box search template to minimize $f : \mathbb{R}^n \to \mathbb{R}$

**Initialize distribution parameters $\theta$, set sample size $\lambda \in \mathbb{N}$**
**While not terminate**

1. **Sample distribution** $P(x|\theta) \to x_1, \ldots, x_\lambda \in \mathbb{R}^n$
2. **Evaluate** $x_1, \ldots, x_\lambda$ **on** $f$
3. **Update parameters** $\theta \leftarrow F_\theta(\theta, x_1, \ldots, x_\lambda, f(x_1), \ldots, f(x_\lambda))$

Everything depends on the definition of $P$ and $F_\theta$

deterministic algorithms are covered as well

In Evolutionary Algorithms the distribution $P$ is often implicitly defined via **operators on a population**, in particular, selection, recombination and mutation

natural template for *Estimation of Distribution Algorithms*

## Metaphors

| Evolutionary Computation | | Optimization |
| --- | --- | --- |
| individual, offspring, parent | $\longleftrightarrow$ | candidate solution |
| | | decision variables |
| | | design variables |
| | | object variables |
| population | $\longleftrightarrow$ | set of candidate solutions |
| fitness function | $\longleftrightarrow$ | objective function |
| | | loss function |
| | | cost function |
| generation | $\longleftrightarrow$ | iteration |

. . . function properties

## The Evolution Strategy

Minimize $f : \mathbb{R}^n \to \mathbb{R}$

**Initialize distribution parameters $\theta$, set population size $\lambda \in \mathbb{N}$**
**While not terminate**

1. **Sample distribution** $P(x|\theta) \to x_1, \ldots, x_\lambda \in \mathbb{R}^n$
2. **Evaluate** $x_1, \ldots, x_\lambda$ **on** $f$
3. **Update parameters** $\theta \leftarrow F_\theta(\theta, x_1, \ldots, x_\lambda, f(x_1), \ldots, f(x_\lambda))$

- $P$ is a **multi-variate normal** distribution

$$\mathcal{N}(m_i, \sigma_i^2 C_i) \sim \boxed{m_i + \sigma_i \mathcal{N}(0, C_i)} \quad \text{for } i = 1, \ldots, \lambda$$

- $\theta = \{m_i, C_i, \sigma_i\}_{i=1,\ldots,\lambda} \in (\mathbb{R}^n \times \mathbb{R}^{n \times n} \times \mathbb{R}_+)^\lambda$
- $F_\theta = F_\theta(\theta, x_{1:\lambda}, \ldots, x_{\mu:\lambda})$, where $\mu \leq \lambda$ and $x_{i:\lambda}$ is the $i$-th best of the $\lambda$ points

...why?

## Why Normal Distributions?

1. widely observed in nature, for example as phenotypic traits

2. only stable distribution with finite variance
   > stable means the *sum* of normal variates is also normal,
   > helpful in **design and analysis** of algorithms

3. most convenient way to generate **isotropic** search points
   > the isotropic distribution does **not favor any direction**
   > (unfoundedly), supports rotational invariance

4. maximum entropy distribution with finite variance
   > the least possible assumptions on $f$ in the distribution shape

## Normal Distribution



Standard Normal Distribution

probability density of 1-D standard normal distribution



2–D Normal Distribution

probability density of 2-D normal distribution

# The Multi-Variate ($n$-Dimensional) Normal Distribution

Any multi-variate normal distribution $\mathcal{N}(\boldsymbol{m}, \mathbf{C})$ is uniquely determined by its mean value $\boldsymbol{m} \in \mathbb{R}^n$ and its symmetric positive definite $n \times n$ covariance matrix $\mathbf{C}$.

The **mean** value $\boldsymbol{m}$

- determines the displacement (translation)
- is the value with the largest density (modal value)
- the distribution is symmetric about the distribution mean



2–D Normal Distribution

The **covariance matrix $\mathbf{C}$**. . .

The **covariance matrix $C$** determines the shape. It has a valuable **geometrical interpretation**: any covariance matrix can be uniquely identified with the iso-density ellipsoid $\{x \in \mathbb{R}^n \mid x^T C^{-1} x = 1\}$

Lines of Equal Density



$\mathcal{N}(m, \sigma^2 I) \sim m + \sigma \mathcal{N}(0, I)$
**one degree of freedom** $\sigma$
components of $\mathcal{N}(0, I)$
are independent standard
normally distributed

$\mathcal{N}(m, D^2) \sim m + D \mathcal{N}(0, I)$
$n$ **degrees of freedom**
components are
independent, scaled

$\mathcal{N}(m, C) \sim m + C^{\frac{1}{2}} \mathcal{N}(0, I)$
$(n^2 + n)/2$ **degrees of freedom**
components are
correlated

where $I$ is the identity matrix (isotropic case) and $D$ is a diagonal matrix (reasonable for separable problems) and $A \times \mathcal{N}(0, I) \sim \mathcal{N}(0, AA^T)$ holds for all $A$.

...CMA

# Sampling New Search Points

The governing equation for derandomized Evolution Strategies

> ### New search points are sampled normally distributed
>
> $$\boldsymbol{x}_i \sim \boldsymbol{m} + \sigma \, \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \qquad \text{for } i = 1, \ldots, \lambda$$
>
> as perturbations of $\boldsymbol{m}$ where $\boldsymbol{x}_i, \boldsymbol{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, and $\mathbf{C} \in \mathbb{R}^{n \times n}$

where

- the mean vector $\boldsymbol{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the *step length*
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

The question remains how to update $\boldsymbol{m}$, $\mathbf{C}$, and $\sigma$.

## Covariance Matrix Adaptation

Rank-One Update

$$\boldsymbol{m} \leftarrow \boldsymbol{m} + \sigma \boldsymbol{y}_w, \quad \boldsymbol{y}_w = \sum_{i=1}^{\mu} w_i \boldsymbol{y}_{i:\lambda}, \quad \boldsymbol{y}_i \sim \mathcal{N}_i(\boldsymbol{0}, \mathbf{C})$$



new distribution,

$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \boldsymbol{y}_w \boldsymbol{y}_w^{\mathrm{T}}$$

the ruling principle: the adaptation increases the probability of successful steps, $\boldsymbol{y}_w$, to appear again

# Rank-$\mu$ Update

$$x_i = m + \sigma y_i, \qquad y_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}),$$
$$m \leftarrow m + \sigma y_w \qquad y_w = \sum_{i=1}^{\mu} w_i y_{i:\lambda}$$



$x_i = m + \sigma y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$

$\mathbf{C}_\mu = \frac{1}{\mu} \sum y_{i:\lambda} y_{i:\lambda}^{\mathsf{T}}$
$\mathbf{C} \leftarrow (1 - 1) \times \mathbf{C} + 1 \times \mathbf{C}_\mu$

$m_{\text{new}} \leftarrow m + \frac{1}{\mu} \sum y_{i:\lambda}$

sampling of $\lambda = 150$
solutions where
$\mathbf{C} = \mathbf{I}$ and $\sigma = 1$

calculating $\mathbf{C}$ from
$\mu = 50$ points,
$w_1 = \cdots = w_\mu = \frac{1}{\mu}$

new distribution

Remark: the old (sample) distribution shape has a great influence on the new
distribution $\longrightarrow$ iterations needed

$$\mathbf{C} \leftarrow (1 - c_{\mathrm{cov}})\mathbf{C} + c_{\mathrm{cov}}\mu_w \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$$

covariance matrix adaptation in the evolution strategy

- learns all **pairwise dependencies** between variables
  off-diagonal entries in the covariance matrix reflect the dependencies

- conducts a **principle component analysis** (PCA) of steps $y_w$,
  sequentially in time and space
  eigenvectors of the covariance matrix $\mathbf{C}$ are the principle components / the
  principle axes of the mutation ellipsoid

- learns a new, **rotated problem represen-
  tation** and a **new metric** (Mahalanobis)
  components are independent (only) in the new representation

- approximates the inverse Hessian on quadratic functions
  overwhelming empirical evidence, proof is in progress

$$\mathbf{C} \leftarrow (1 - c_{\mathrm{cov}})\mathbf{C} + c_{\mathrm{cov}}\mu_w \mathbf{y}_w \mathbf{y}_w^{\mathrm{T}}$$

covariance matrix adaptation

- is equivalent with an adaptive (general) linear
  encoding[1]

---

[1] Hansen 2000, Invariance, Self-Adaptation and Correlated Mutations in Evolution Strategies, PPSN VI

**Problem Statement**
000

**The Challenges**
000

**Evolution Strategy**
00000000

**CMA**
0000●0

**Evaluation**
00000

**Adaptive Encoding**
0000000

# Experimentum Crucis (1)

$f$ convex quadratic, separable



$$f(\boldsymbol{x}) = \sum_{i=1}^{n} 10^{\alpha \frac{i-1}{n-1}} x_i^2, \alpha = 6$$

## Experimentum Crucis (2)

$f$ convex quadratic, as before but non-separable (rotated)



$\mathbf{C} \propto \boldsymbol{H}^{-1}$ **for all** $g, \mathbf{H}$

$f(\boldsymbol{x}) = g\left(\boldsymbol{x}^{\mathrm{T}} \mathbf{H} \boldsymbol{x}\right), \, g : \mathbb{R} \to \mathbb{R}$ stricly monotonic

## Invariance

> *The grand aim of all science is to cover the greatest number of empirical facts by logical deduction from the smallest number of hypotheses or axioms.*
>
> — Albert Einstein

**Invariance** is a *guaranty for generalization* of performance from a single function to a class of functions. Most important invariance properties of the Covariance Matrix Adaptation (CMA) Evolution Strategy (ES) are

- invariance to **order preserving transformations** in function space $\longrightarrow$
- Translation and **rotation invariance** in search space
  to *rigid transformations* of the search space



... empirical validation

# Comparison to , BFGS, PSO and DE

$f$ convex quadratic, non-separable (rotated) with varying $\alpha$



Ellipsoid dimension 20, 21 trials, tolerance 1e−09, eval max 1e+07

SP1 = average number of objective function evaluations to reach the target function value of $10^{-9}$

$f(\boldsymbol{x}) = g(\boldsymbol{x}^{\mathrm{T}} \mathbf{H} \boldsymbol{x})$ with

$g$ identity (BFGS, *red*) or

$g(.) = (.)^{1/4}$ (BFGS, *red dashed*) or

$g$ order-preserving = strictly increasing (all other)

BFGS: quasi-Newton method
PSO: Particle Swarm Optimization
DE:Differential Evolution
CMA-ES —

...population size, invariance

# Comparison to IDEA and Simplex-Downhill



CMA-ES: Covariance Matrix Adaptation Evolution Strategy[2]
IDEA: Iterated Density-Estimation Evolutionary Algorithm[3]
Fminsearch: Nelder-Mead simplex downhill method[4]
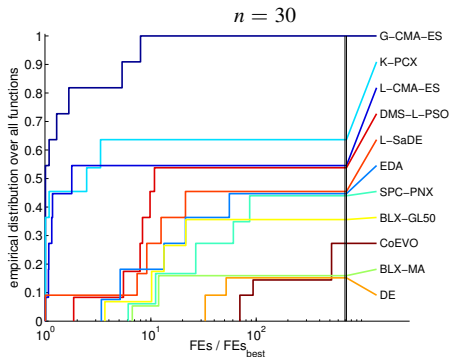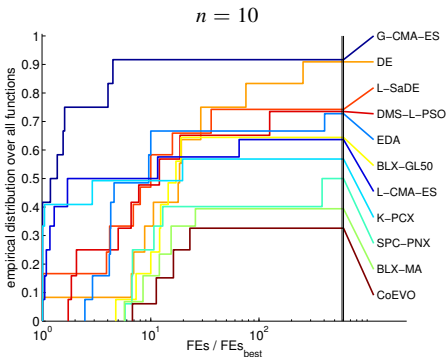
Randomsearch: pure Monte-Carlo sampling

Peter Dürr and Andreas Pfister 2004. Optimization of Walking Gaits for a Three Legged Walking Robot,

Diploma Thesis, Institut für Mechanische Systeme, ETH Zurich

---

[2] Hansen (2001) Completely Derandomized Self-Adaptation in Evolution Strategies. *Evolutionary Computation Journal*

[3]

# CEC 2005 Real Parameter Optimization Session

Empirical Distribution of Normalized Success Performance



$$\mathrm{FEs} = \mathrm{mean}(\#\textit{fevals}) \times \frac{\#\text{all runs (25)}}{\#\text{successful runs}}, \text{ where } \#\textit{fevals} \text{ includes only successful runs.}$$

Shown: **empirical distribution function** of the Success Performance FEs divided by FEs of the best algorithm on the respective function.

Results of all functions are used where at least one algorithm was successful at least once, i.e. where the target function value was reached in at least one experiment (out of $11 \times 25$ experiments).

Small values for FEs and therefore large (cumulative frequency) values in the graphs are preferable.

# The Covariance Matrix Adaptation Evolution Strategy

In a Nutshell

1. Multivariate normal distribution to generate new search points

   *follows the maximum entropy principle*

2. Selection only based on the ranking of the $f$-values, weighted recombination

   *using only the ranking of $f$-values preserves invariance*

3. *Covariance matrix adaptation (CMA)* **increases the probability** to repeat **successful steps**

   *learning all pairwise dependencies*
   $\implies$ *conducts an incremental PCA*
   $\implies$ *new (rotated) problem representation*

4. An **evolution path** (a trajectory) is exploited in two places
   - enhances the covariance matrix (rank-one) adaptation

     *yields sometimes linear time complexity*
   - controls the **step-size** (step length)

     *aims at conjugate perpendicularity*

## Linear Encoding and the Covariance Matrix

Equivalence between change in encoding and transformation of the mutation operator

Let $x_B, x_A \in \mathbb{R}^n$ be **two genotypes** encoding the same phenotype

$$y = A\,x_A = B\,x_B$$

The **effect of the different encodings** becomes evident, when the genotype is changed (adding $\mathcal{N}(\mathbf{0}, \mathbf{C})$).

$$
\begin{aligned}
y_{\text{new}} &= B\,(x_B + \mathcal{N}(\mathbf{0}, \mathbf{C})) &= B\,x_B + B\,\mathcal{N}(\mathbf{0}, \mathbf{C}) \\
&&= A\,x_A + A\,A^{-1}B\,\mathcal{N}(\mathbf{0}, \mathbf{C}) \\
&&= A\,(x_A + A^{-1}B\,\mathcal{N}(\mathbf{0}, \mathbf{C})) \\
y_{\text{new}} &= A\,(x_A + A^{-1}B\,\mathcal{N}(\mathbf{0}, \mathbf{C}))
\end{aligned}
$$

Using a new encoding $B$ means using a different covariance matrix

## Adaptive Encoding

### Definition (Adaptive Encoding)

Given a search algorithm, $\mathcal{A}$ in state $s$, an encoding, $T_B$ and an update, $\mathcal{U}$, then the iteration step

$$s \quad \leftarrow \quad T_B \circ \mathcal{A} \circ T_B^{-1}(s) \tag{1}$$
$$\mathbf{B} \quad \leftarrow \quad \mathcal{U}(\mathbf{B}, s) \tag{2}$$

defines an *adaptive encoding* where $T_B \circ \mathcal{A} \circ T_B^{-1}(s) = T_B(\mathcal{A}(T_B^{-1}(s)))$.

### Remark (Evaluation of Solutions)

In order to make use of Eq. (1), $\mathcal{A}$ has to operate on $f \circ \mathbf{B}$.

# Adaptive Encoding

Example: Adaptive Encoding of CSA-ES

## $AE_{CMA}$-CSA-ES

1 initialize $m \in \mathbb{R}^n$ (distribution mean), $p_\sigma = 0$ (evolution path), $\sigma > 0$ (step-size)

2 initialize $B = B^\circ = I$ (encoding matrices)

3 **repeat**

4     $m \leftarrow B^{-1} m$

5     $p_\sigma \leftarrow B^{\circ \mathrm{T}} p_\sigma$

6     **begin**

7        $x_i = m + \sigma \mathcal{N}_i(0, I)$,   for $i = 1, \ldots, \lambda$

8        $f_i = f \circ B(x_i) = f(B x_i)$, for $i = 1, \ldots, \lambda$   // encode to evaluate

9        $m^- = m$

10       $m \leftarrow \sum_{i=1}^\mu w_i x_{i(f)}$

11       $p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \sqrt{c_\sigma (2 - c_\sigma) \mu_w} \frac{1}{\sigma}(m - m^-)$

12       $\sigma \leftarrow \sigma \exp\left( \frac{c_\sigma}{d_\sigma} \left( \frac{\|p_\sigma\|}{E \|\mathcal{N}(0, I)\|} - 1 \right) \right)$

13     **end**

14     $m \leftarrow B m$

15     $p_\sigma \leftarrow B^\circ p_\sigma$

16     $AE_{CMA}$-Update($\{B x_1, \ldots, B x_\mu\}$)            // update $B$ and $B^\circ$

17 **until** *stopping criterion is met*

        AE : Adaptive Encoding

      CMA : Covariance Matrix Adaptation

CSA-ES : Cumulative *Step-size Adaptation* Evolution Strategy, lines 6–13

## Adaptive Encoding

### Theorem (Recovery of CMA-ES)

*Given $AE_{CMA}$-Update in Procedure 1, the $AE_{CMA}$-$(\mu/\mu_W, \lambda)$-CSA-ES implements the $(\mu/\mu_W, \lambda)$-CMA-ES.*

Adaptive Encoding

- can render *any* continuous domain search algorithm independent of the coordinate system

- anticipated successful applications in particular for population-based stochastic algorithms
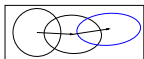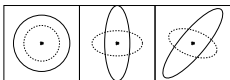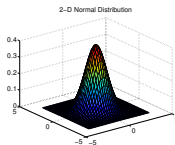
## Another Case Study

Adaptive Encoding



Cauchy-ES (green) versus Adaptively Encoded Cauchy-ES (black)

rotating a separable function, 10-D

## Another Case Study

Adaptive Encoding



Cauchy-ES (green) versus Adaptively Encoded Cauchy-ES (black)

rotating a separable function, 30-D

# Thank You