

The difficulties of black-box optimization and a stochastic variable metric approach

Nikolaus Hansen

INRIA – Saclay, TAO research group

November 9, 2009

slides: <http://www.lri.fr/~hansen/ParisRoc-handout.pdf>

Content

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
- 5 Convergence Properties
- 6 Performance Evaluation

*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

Continuous Domain Search/Optimization

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

$$f : \mathcal{X} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}, \quad \underline{x} \mapsto f(\underline{x})$$

- **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

Continuous Domain Search/Optimization

- Goal

- ▶ solution \underline{x} with **small function value** with **least search cost**
there are two conflicting objectives
- ▶ fast convergence to the global optimum
... or to a robust solution \underline{x}

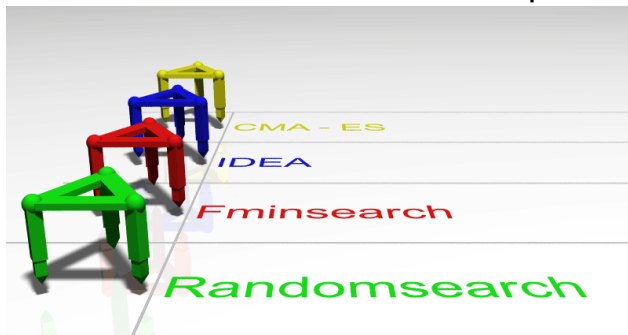
- Typical Examples

- ▶ shape optimization (e.g. using CFD) curve fitting, airfoils
- ▶ model calibration biological, physical
- ▶ parameter calibration algorithms, controllers,
plants, images

Approach: stochastic search, Evolutionary Algorithms

... tripods

Comparison of CMA-ES, IDEA and Simplex-Downhill



CMA-ES: Covariance Matrix Adaptation Evolution Strategy

IDEA: Iterated Density-Estimation Evolutionary Algorithm¹

Fminsearch: Nelder-Mead simplex downhill method²

P. Dürr and A. Pfister (2004), Optimization of Walking Gaits for a Three Legged Robot, term paper.

see...<http://www.icos.ethz.ch/cse/research/highlights/research.highlights-august-2004>

¹ Bosman (2003) Design and Application of Iterated Density-Estimation Evolutionary Algorithms. PhD thesis.

² Nelder and Mead (1965). A simplex method for function minimization. *Computer Journal*.

1 Problem Statement

2 The Difficulties

3 Stochastic Search

4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)

- Sampling
- Step-Size Control
- Covariance Matrix Adaptation
- Cumulation—the Evolution Path
- Covariance Matrix Rank- μ Update

5 Convergence Properties

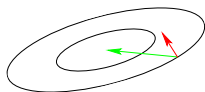
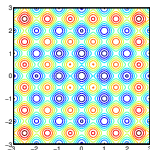
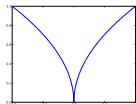
6 Performance Evaluation

... metaphores

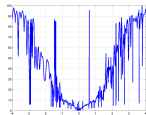
What Makes a Function Difficult to Solve?

Why stochastic search?

- non-linear, non-quadratic, non-convex
on linear/quadratic functions better search policies are available
- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning
widely varying sensitivity
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function



gradient direction Newton direction



Curse of Dimensionality

The term *Curse of dimensionality* (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding dimensions to a (mathematical) space.

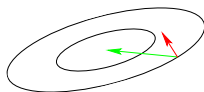
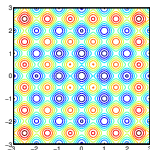
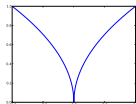
Example: Consider placing 100 points onto a real interval, say $[0, 1]$. To get **similar coverage**, in terms of distance between adjacent points, of the 10-dimensional space $[0, 1]^{10}$ would require $100^{10} = 10^{20}$ points. A 100 points have minimal distance of ≈ 0.65 (on average) and appear now as isolated points in a vast empty space.

Implication: A **search policy** (e.g. exhaustive search) that is **efficient** in small dimensions **might be useless** in moderate or large dimensional search spaces.

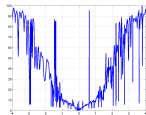
What Makes a Function Difficult to Solve?

Why stochastic search?

- non-linear, non-quadratic, non-convex
on linear/quadratic functions better search policies are available
- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning
widely varying sensitivity
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function



gradient direction Newton direction



Separable Problems

Definition (Separable Problem)

A function f is separable if

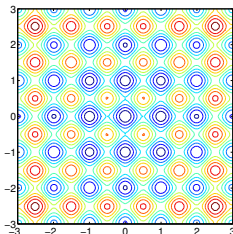
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left(\arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

⇒ it follows that f can be optimized in a sequence of n independent 1-D optimization processes

Example: Additively decomposable functions

$$f(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i)$$

Rastrigin function



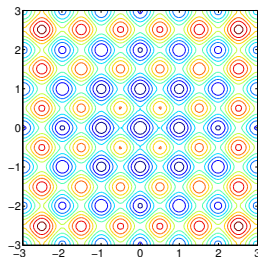
Non-Separable Problems

Building a non-separable problem from a separable one

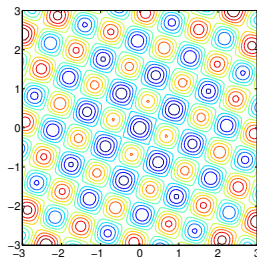
Rotating the coordinate system

- $f : \underline{x} \mapsto f(\underline{x})$ separable
- $f : \underline{x} \mapsto f(\underline{Rx})$ **non-separable**

R rotation matrix



R
→



34

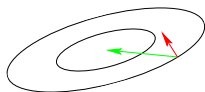
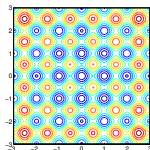
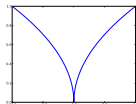
³ Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

⁴ Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

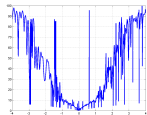
What Makes a Function Difficult to Solve?

Why stochastic search?

- non-linear, non-quadratic, non-convex
on linear/quadratic functions better search policies are available
- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning
widely varying sensitivity
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function



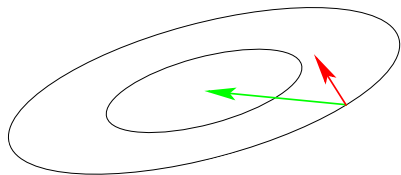
gradient direction Newton direction



III-Conditioned Problems

Curvature of level sets

Consider the convex-quadratic function $f(\underline{x}) = \frac{1}{2}(\underline{x} - \underline{x}^*)^T \underline{\underline{H}}(\underline{x} - \underline{x}^*)$



gradient direction $-f'(\underline{x})^T$

Newton direction $-\underline{\underline{H}}^{-1}f'(\underline{x})^T$

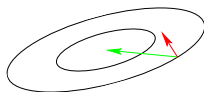
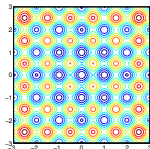
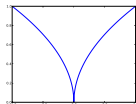
Condition number equals nine here. Condition numbers between 100 and even 10^{10} can often be observed in real world problems.

If $\underline{\underline{H}} \approx \underline{\underline{I}}$ (small condition number of $\underline{\underline{H}}$) first order information (e.g. the gradient) is sufficient. Otherwise **second order information** (estimation of $\underline{\underline{H}}^{-1}$) **is required**.

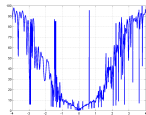
What Makes a Function Difficult to Solve?

Why stochastic search?

- non-linear, non-quadratic, non-convex
on linear/quadratic functions better search policies are available
- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning
widely varying sensitivity
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function

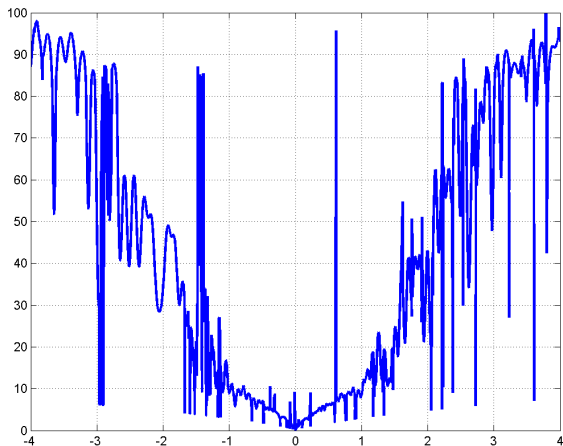


gradient direction Newton direction



Ruggedness

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



cut from an (easily) solvable example in 5-D

What Makes a Function Difficult to Solve?

... and what can be done

Challenge	Approach in Evolutionary Computation
Dimensionality, Non-Separability	exploiting the problem structure locality, neighborhood, encoding
Ill-conditioning	second order approach changes the neighborhood metric
Ruggedness	non-local policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed stochastic, non-elitistic, population-based method recombination operator serves as repair mechanism

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search**
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

... metaphores

Metaphors

(Biological) **Evolution**(ary Computation)

Optimization

genome

↔

decision variables

design variables

object variables

individual, offspring, parent

↔

candidate solution

population

↔

set of candidate solutions

fitness function

↔

objective function

loss function

cost function

generation

↔

iteration

Stochastic Search

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

Initialize distribution parameters $\underline{\theta}$, set population size $\lambda \in \mathbb{N}$

While not terminate

- ① **Sample distribution** $P(\underline{x}|\underline{\theta}) \rightarrow \underline{x}_1, \dots, \underline{x}_\lambda \in \mathbb{R}^n$
- ② **Evaluate** $\underline{x}_1, \dots, \underline{x}_\lambda$ **on** f
- ③ **Update parameters** $\underline{\theta} \leftarrow F_\theta(\underline{\theta}, \underline{x}_1, \dots, \underline{x}_\lambda, f(\underline{x}_1), \dots, f(\underline{x}_\lambda))$

Everything depends on the definition of P and F_θ

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*

Stochastic Search

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

Initialize distribution parameters $\underline{\theta}$, set population size $\lambda \in \mathbb{N}$

While not terminate

- ① **Sample distribution** $P(\underline{x}|\underline{\theta}) \rightarrow \underline{x}_1, \dots, \underline{x}_\lambda \in \mathbb{R}^n$
- ② **Evaluate** $\underline{x}_1, \dots, \underline{x}_\lambda$ on f
- ③ **Update parameters** $\underline{\theta} \leftarrow F_\theta(\underline{\theta}, \underline{x}_1, \dots, \underline{x}_\lambda, f(\underline{x}_1), \dots, f(\underline{x}_\lambda))$

In the following

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

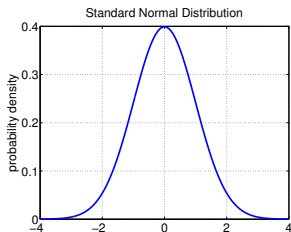
$$\mathcal{N}(\underline{m}, \sigma^2 \underline{C}) \sim \underline{m} + \sigma \mathcal{N}(\underline{0}, \underline{C})$$

$$\underline{\theta} = \{\underline{m}, \underline{C}, \sigma\} \in \mathbb{R}^n \times \mathbb{R}^{n \times n} \times \mathbb{R}_+$$

- $F_\theta = F_\theta(\underline{\theta}, \underline{x}_{1:\lambda}, \dots, \underline{x}_{\mu:\lambda})$, where $\mu \leq \lambda$ and $\underline{x}_{i:\lambda}$ is the i -th best of the λ points

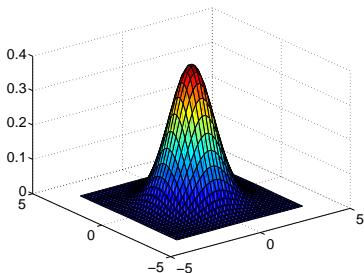
... normal distribution

Normal Distribution



probability density of 1-D standard normal distribution

2-D Normal Distribution



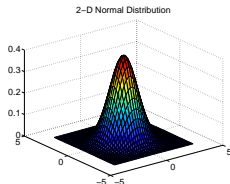
probability density of 2-D normal distribution

The Multi-Variate (n -Dimensional) Normal Distribution

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

The **mean** value \underline{m}

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

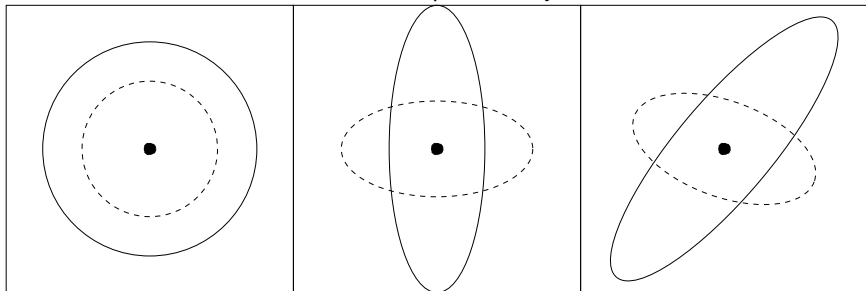


The **covariance matrix** $\underline{\underline{C}}$

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

... any **covariance matrix** can be uniquely identified with the iso-density ellipsoid
 $\{\underline{x} \in \mathbb{R}^n \mid \underline{x}^T \underline{C}^{-1} \underline{x} = 1\}$

Lines of Equal Density



$$\mathcal{N}(\underline{m}, \sigma^2 \underline{I}) \sim \underline{m} + \sigma \mathcal{N}(\underline{0}, \underline{I})$$

one degree of freedom σ
 components are
 independent standard
 normally distributed

$$\mathcal{N}(\underline{m}, \underline{D}^2) \sim \underline{m} + \underline{D} \mathcal{N}(\underline{0}, \underline{I})$$

n degrees of freedom
 components are
 independent, scaled

$$\mathcal{N}(\underline{m}, \underline{C}) \sim \underline{m} + \underline{C}^{\frac{1}{2}} \mathcal{N}(\underline{0}, \underline{I})$$

$(n^2 + n)/2$ degrees of freedom
 components are
 correlated

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

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

Rank-Based Stochastic Search

Rank-based black box search to minimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters $\underline{\theta}$, set population size $\lambda \in \mathbb{N}$

While not terminate

① **Sample distribution** $P(\underline{x}|\underline{\theta}) \rightarrow \underline{x}_1, \dots, \underline{x}_\lambda \in \mathbb{R}^n$

② **Evaluate** $\underline{x}_1, \dots, \underline{x}_\lambda$ on f

and let $f(\underline{x}_{i:\lambda}) \leq f(\underline{x}_{j:\lambda}) \Leftrightarrow i \leq j$

③ **Update parameters** $\underline{\theta} \leftarrow F_\theta(\underline{\theta}, \underline{x}_{1:\lambda}, \dots, \underline{x}_{\mu:\lambda})$

P is a multi-variate normal distribution $\mathcal{N}\left(\underline{m}, \sigma^2 \underline{C}\right) \sim \underline{m} + \sigma \mathcal{N}\left(\underline{0}, \underline{C}\right)$

... sampling

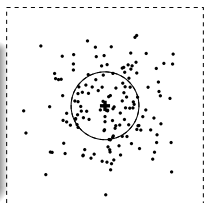
Sampling New Search Points

The Mutation Operator

New search points are normally distributed

$$\underline{x}_i \sim \underline{m} + \sigma \mathcal{N}_i(\underline{0}, \underline{C}) \quad \text{for } i = 1, \dots, \lambda$$

perturbations of \underline{m}



where

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

The question remains how to update \underline{m} , σ , and \underline{C} .

Update of the Distribution Mean \underline{m}

Selection and Recombination

The best μ points are selected from the sampled solutions (non-elitistic) and a **weighted mean** is taken:

Given the i -th solution point $\underline{x}_i = \underline{m} + \sigma \underbrace{\mathcal{N}_i(\underline{0}, \underline{C})}_{=: \underline{y}_i} = \underline{m} + \sigma \underline{y}_i$

Let $\underline{x}_{i:\lambda}$ the **i -th ranked** solution point, such that $f(\underline{x}_{1:\lambda}) \leq \dots \leq f(\underline{x}_{\lambda:\lambda})$.

The new mean reads

$$\underline{m} \leftarrow \sum_{i=1}^{\mu} w_i \underline{x}_{i:\lambda} = \underline{m} + \sigma \underbrace{\sum_{i=1}^{\mu} w_i \underline{y}_{i:\lambda}}_{=: \underline{y}_w}$$

where

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

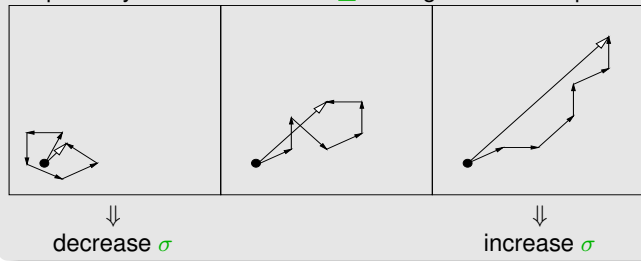
- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

Path Length Control

The Concept

Measure the length of the *evolution path*

the pathway of the mean vector \underline{m} in the generation sequence



$$\underline{x}_i = \underline{m} + \sigma \underline{y}_{-i}$$

$$\underline{m} \leftarrow \underline{m} + \sigma \underline{y}_{-w}$$

loosely speaking steps are

- perpendicular under random selection (in expectation)
- perpendicular in the desired situation (to be most efficient)

Path Length Control: Equations

AKA Cumulative Step-Size Adaptation (CSA)

Initialize $\underline{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\underline{C} = \underline{I}$, and $\underline{p}_\sigma = \underline{0}$

set $c_\sigma \approx 4/n$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$,

set λ and $w_{i=1,\dots,\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^\lambda w_i^2} \approx 0.3 \lambda$

While not terminate

$\underline{x}_i = \underline{m} + \sigma \underline{y}_i$, where $\underline{y}_i \sim \mathcal{N}(\underline{0}, \underline{C})$ for $i = 1, \dots, \lambda$ sampling

$\underline{m} \leftarrow \underline{m} + \sigma \underline{y}_w$ where $\underline{y}_w = \sum_{i=1}^\mu w_i \underline{y}_{i:\lambda}$ update mean

$\underline{p}_\sigma \leftarrow (1 - c_\sigma) \underline{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} \underline{C}^{-\frac{1}{2}} \underline{y}_w$ cumulation for σ

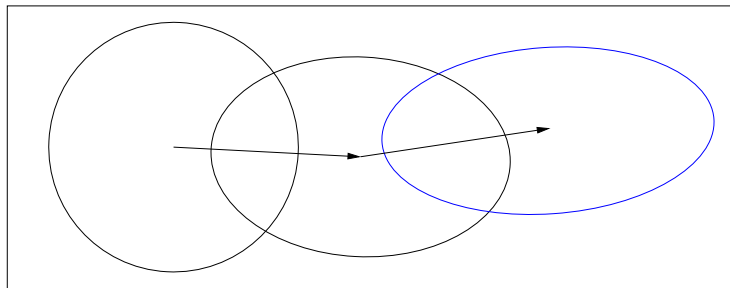
$\sigma \leftarrow \sigma \times \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\underline{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\underline{0}, \underline{I})\|} - 1\right)\right)$ update of σ

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

Covariance Matrix Adaptation

Rank-One Update

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



new distribution,

$$\underline{C} \leftarrow 0.8 \times \underline{C} + 0.2 \times \underline{y}_{-w} \underline{y}_{-w}^T$$

the ruling principle: the adaptation **increases the likelihood of successful steps**, \underline{y}_{-w} , to appear again

... equations

Preliminary Set of Equations

Covariance Matrix Adaptation with Rank-One Update

Initialize $\underline{m} \in \mathbb{R}^n$, and $\underline{C} = \underline{I}$, set $\sigma = 1$, learning rate $c_{\text{cov}} \approx 2/n^2$

While not terminate

$$\underline{x}_i = \underline{m} + \sigma \underline{y}_i, \quad \underline{y}_i \sim \mathcal{N}_i(\underline{0}, \underline{C}), \quad i = 1, \dots, \lambda$$

$$\underline{m} \leftarrow \underline{m} + \sigma \underline{y}_{-w} \quad \text{where } \underline{y}_{-w} = \sum_{i=1}^{\mu} w_i \underline{y}_{i:\lambda}$$

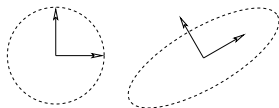
$$\underline{C} \leftarrow (1 - c_{\text{cov}}) \underline{C} + c_{\text{cov}} \underbrace{\mu_w \underline{y}_{-w} \underline{y}_{-w}^T}_{\text{rank-one}} \quad \text{where } \mu_w = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \geq 1$$

λ can be small

$$\underline{\underline{C}} \leftarrow (1 - c_{\text{cov}})\underline{\underline{C}} + c_{\text{cov}}\mu_w y_{-w} y_{-w}^T$$

The covariance matrix adaptation

- 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 $\underline{\underline{C}}$ are the principle components / the principle axes of the mutation ellipsoid
- learns a new, **rotated problem representation** and a **new variable metric** (Mahalanobis)
components are independent (only) in the new representation
rotational invariant
equivalent with an adaptive (general) linear encoding
- approximates the **inverse Hessian** on convex-quadratic functions
overwhelming empirical evidence, proof is in progress

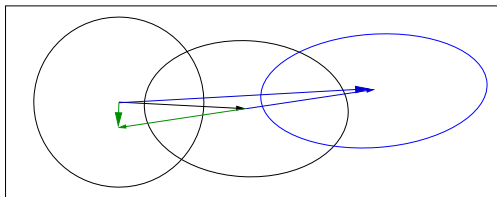


- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

Cumulation

Utilizing the Evolution Path

We used $\underline{y}_{-w} \underline{y}_{-w}^T$ for updating \underline{C} . Because $\underline{y}_{-w} \underline{y}_{-w}^T = -\underline{y}_{-w} (-\underline{y}_{-w})^T$ **the sign** of \underline{y}_{-w} is neglected. The sign information is (re-)introduced by using the *evolution path*.



$$\underline{p}_{-c} \leftarrow \underbrace{(1 - c_c)}_{\text{decay factor}} \underline{p}_{-c} + \underbrace{\sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w}}_{\text{normalization factor}} \underline{y}_{-w}$$

$$\underline{C} \leftarrow (1 - c_{\text{cov}}) \underline{C} + c_{\text{cov}} \underbrace{\underline{p}_{-c} \underline{p}_{-c}^T}_{\text{rank-one}}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_c \ll 1$.

... $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$

Using an **evolution path** for the **rank-one update** of the covariance matrix reduces the number of function evaluations to adapt to a straight ridge **from** $\mathcal{O}(n^2)$ **to** $\mathcal{O}(n)$.^a

^aHansen, Müller and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18

The overall model complexity is n^2 but important parts of the model can be learned in time of order n

... rank- μ

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - Sampling
 - Step-Size Control
 - Covariance Matrix Adaptation
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
- 5 Convergence Properties
- 6 Performance Evaluation

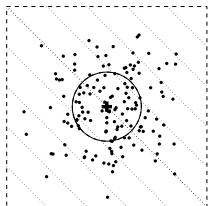
Rank- μ Update

$$\begin{aligned} \underline{x}_i &= \underline{m} + \sigma \underline{y}_i, & \underline{y}_i &\sim \mathcal{N}_i(\underline{0}, \underline{C}), \\ \underline{m} &\leftarrow \underline{m} + \sigma \underline{y}_{-w}, & \underline{y}_{-w} &= \sum_{i=1}^{\mu} w_i \underline{y}_{i:\lambda} \end{aligned}$$

The rank- μ update extends the update rule for **large population sizes** λ using $\mu > 1$ vectors to update \underline{C} at each generation step.

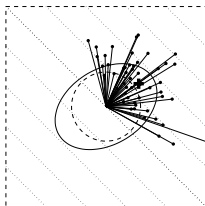
$$\underline{C} \leftarrow (1 - c_{\text{cov}}) \underline{C} + c_{\text{cov}} \sum_{i=1}^{\mu} w_i \underline{y}_{i:\lambda} \underline{y}_{i:\lambda}^T$$

where $c_{\text{cov}} \approx \mu_w/n^2$ and $c_{\text{cov}} \leq 1$.



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

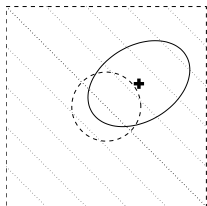
sampling of $\lambda = 150$
solutions where
 $\underline{C} = \underline{I}$ and $\sigma = 1$



$$\underline{C}_{\mu} = \frac{1}{\mu} \sum y_{i:\lambda} y_{i:\lambda}^T$$

$$\underline{C} \leftarrow (1 - 1/\mu) \times \underline{C} + 1/\mu \times \underline{C}_{\mu}$$

calculating \underline{C} where
 $\mu = 50$,
 $w_1 = \dots = w_{\mu} = \frac{1}{\mu}$,
and $c_{\text{cov}} = 1$



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

new distribution

The rank- μ update

- increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- can reduce the number of necessary **generations** roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁵
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

⁵Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18

Summary of Equations

The Covariance Matrix Adaptation Evolution Strategy

Input: $\underline{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, λ

Initialize $\underline{C} = \underline{I}$, and $\underline{p}_c = \underline{0}$, $\underline{p}_\sigma = \underline{0}$,

set $c_c \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n^2$, $c_\mu \approx \mu_w/n^2$, $c_1 + c_\mu \leq 1$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$, and $w_{i=1,\dots,\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^\mu w_i^2} \approx 0.3 \lambda$

While not terminate

$$\underline{x}_i = \underline{m} + \sigma \underline{y}_i, \quad \underline{y}_i \sim \mathcal{N}_i(\underline{0}, \underline{C}), \quad \text{for } i = 1, \dots, \lambda \quad \text{sampling}$$

$$\underline{m} \leftarrow \underline{m} + \sigma \underline{y}_w \quad \text{where } \underline{y}_w = \sum_{i=1}^\mu w_i \underline{y}_{i:\lambda} \quad \text{update mean}$$

$$\underline{p}_c \leftarrow (1 - c_c) \underline{p}_c + \mathbf{1}_{\{\|\underline{p}_c\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \underline{y}_w \quad \text{cumulation for } \underline{C}$$

$$\underline{p}_\sigma \leftarrow (1 - c_\sigma) \underline{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} \underline{C}^{-\frac{1}{2}} \underline{y}_w \quad \text{cumulation for } \sigma$$

$$\underline{C} \leftarrow (1 - c_1 - c_\mu) \underline{C} + c_1 \underline{p}_c \underline{p}_c^T + c_\mu \sum_{i=1}^\mu w_i \underline{y}_{i:\lambda} \underline{y}_{i:\lambda}^T \quad \text{update } \underline{C}$$

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

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
- 5 Convergence Properties**
- 6 Performance Evaluation

Experimentum Crucis

What did we want to achieve?

- reduce any convex-quadratic function

$$f(\underline{x}) = \underline{x}^T \underline{H} \underline{x}$$

e.g. $f(\underline{x}) = \sum_{i=1}^n 10^{6 \frac{i-1}{n-1}} x_i^2$

to the sphere model

$$f(\underline{x}) = \underline{x}^T \underline{x}$$

without use of derivatives

- lines of equal density align with lines of equal fitness

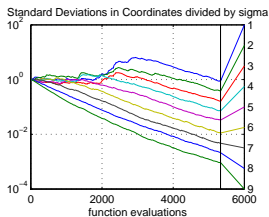
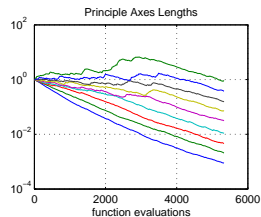
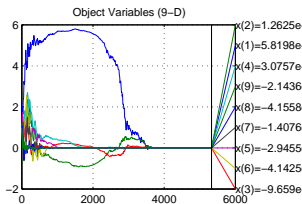
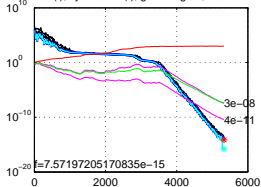
$$\underline{C} \propto \underline{H}^{-1}$$

in a stochastic sense

Experimentum Crucis (1)

f convex-quadratic, separable

blue:abs(f), cyan:f-min(f), green:sigma, red:axis ratio

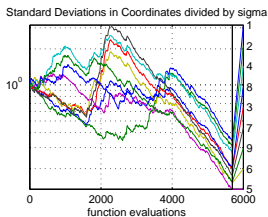
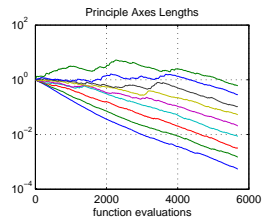
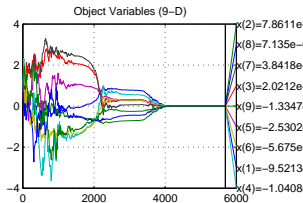
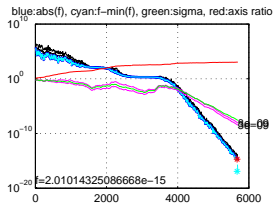


$$f(\underline{x}) = \sum_{i=1}^n 10^{\alpha \frac{i-1}{n-1}} x_i^2, \alpha = 6$$

... crucis rotated

Experimentum Crucis (2)

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



$$\underline{C} \propto \underline{H}^{-1} \text{ for all } g, \underline{H}$$

$$f(\underline{x}) = g(\underline{x}^T \underline{H} \underline{x}), \quad g: \mathbb{R} \rightarrow \mathbb{R} \text{ strictly monotonic}$$

... on convergence

On Convergence

- convergence on a very **broad class of functions**, e.g. for Monte Carlo pure random search

$$\text{very slow: } \|\underline{x} - \underline{x}^*\| \propto \frac{1}{t^{1/n}} = \frac{1}{\exp\left(\frac{\log t}{n}\right)}$$

- convergence with practically **feasible convergence rates** on, e.g., $g\left(\frac{1}{2}\underline{x}^T \underline{\underline{H}} \underline{x}\right)$

$$\text{CMA-ES} \Rightarrow \|\underline{x} - \underline{x}^*\| \propto \frac{1}{\exp\left(\frac{t/4}{n}\right)}$$

- 1 Problem Statement
- 2 The Difficulties
- 3 Stochastic Search
- 4 The Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
- 5 Convergence Properties
- 6 Performance Evaluation**

Evaluation/Selection of Search Algorithms

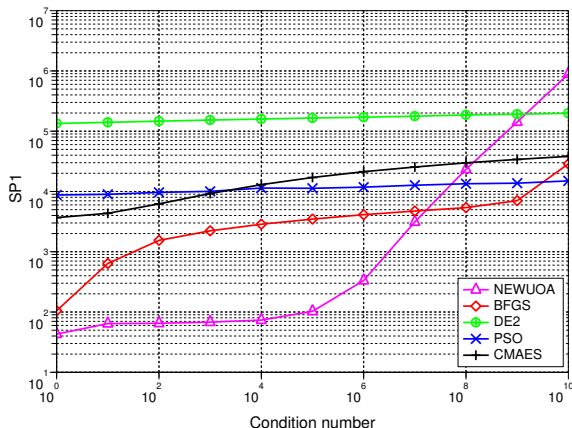
Evaluation (of the performance) of a search algorithm needs

- meaningful **quantitative measure** on benchmark functions or real world problems
- account for **meta-parameter tuning**
can be quite expensive
- account for **invariance properties** (symmetries)
prediction of performance is based on “similarity”, ideally equivalence classes of functions
- account for **algorithm internal cost**
often negligible, depending on the objective function cost

Comparison to BFGS, NEWUOA, PSO and DE (1)

f convex-quadratic, separable with varying α

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



$f(\underline{x}) = g(\underline{x}^T \underline{H} \underline{x})$ with
 g identity (BFGS, NEWUOA) or
 g order-preserving (strictly increasing, all other)

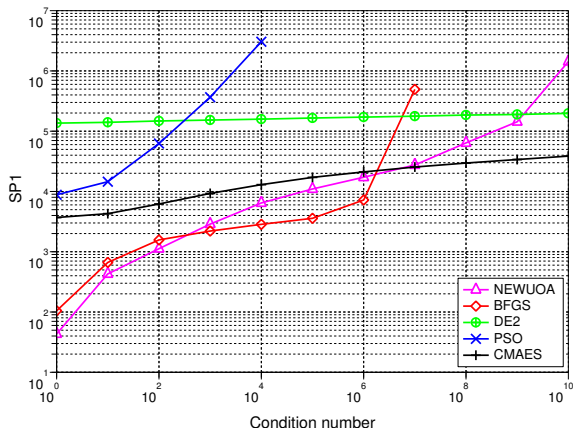
SP1 = average number of objective function evaluations to reach the target function value of 10^{-9}

... population size, invariance

Comparison to BFGS, NEWUOA, PSO and DE (2)

f convex-quadratic, non-separable (rotated) with varying α

Rotated Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



$f(\underline{x}) = g(\underline{x}^T \underline{H} \underline{x})$ with
 g identity (BFGS, NEWUOA) or
 g order-preserving (strictly increasing, all other)

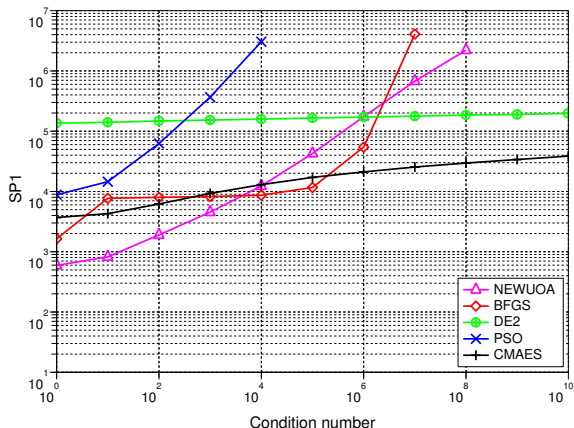
SP1 = average number of objective function evaluations to reach the target function value of 10^{-9}

... population size, invariance

Comparison to BFGS, NEWUOA, PSO and DE (3)

f non-convex, non-separable (rotated) with varying α

Sqrt of sqrt of rotated ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



$f(x) = g(x^T \underline{H} x)$ with
 $g(\cdot) = (\cdot)^{1/4}$ (BFGS, NEWUOA) or
 g any order-preserving (strictly increasing, all other)

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

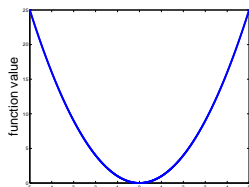
... population size, invariance

Invariance

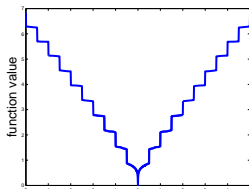
The short version

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.

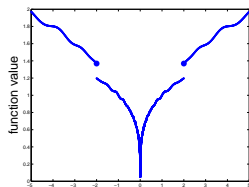
— Albert Einstein



$$f(x) = x^2$$



$$f(x) = g_1(x^2)$$

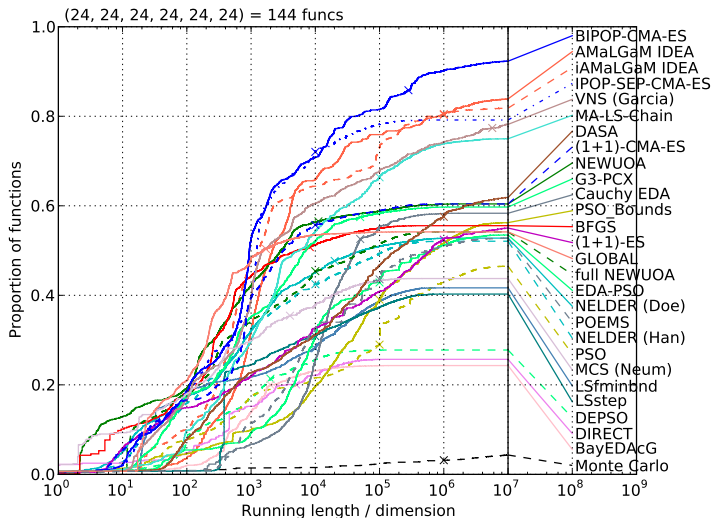


$$f(x) = g_2(x^2)$$

- all three functions are **equivalent** for rank-based search methods
large equivalence class
- invariance allows a save **generalization** of empirical results
here on $f(x) = x^2$ (left) to any $f(x) = g(x^2)$, where g is monotonous

Comprehensive Comparison of 28 Algorithms

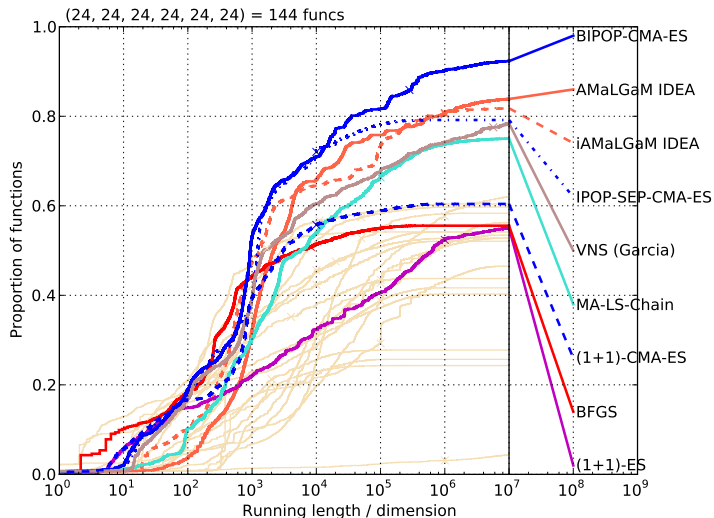
Empirical Distribution of Expected Running Length



on 24 benchmark functions of BBOB 2009 in 20-D

Comprehensive Comparison of 28 Algorithms

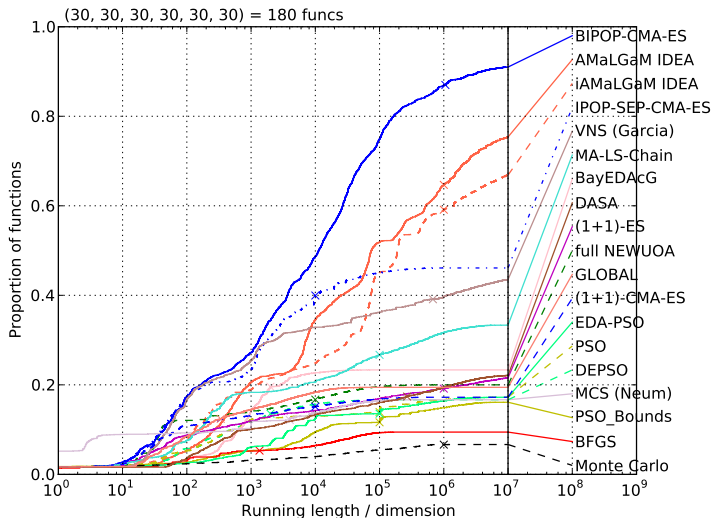
Empirical Distribution of Expected Running Length



on 24 benchmark functions of BBOB 2009 in 20-D

Comprehensive Comparison of 19 Algorithms

Empirical Distribution of Expected Running Length



on 30 noisy benchmark functions of BBOB 2009 in 20-D

Summary

Covariance Matrix Adaptation Evolution Strategy (CMA-ES) 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
preserves invariance
- ③ *Covariance matrix adaptation (CMA)* **increases the likelihood of previously successful steps**
learning all pairwise dependencies
⇒ adapts a variable metric
⇒ new (rotated) problem representation
- ④ An **evolution path** (a non-local trajectory)
 - ▶ enhances the covariance matrix (rank-one) adaptation
yields sometimes linear time complexity
 - ▶ controls the **step-size** (step length)
aims at conjugate perpendicularity

Merci !

<http://www.lri.fr/~hansen/cmaesintro.html>
or google NIKOLAUS HANSEN