

Introduction to Optimization

Lecture 5: CMA-ES

October 6, 2017
TC2 - Optimisation
Université Paris-Saclay



Dimo Brockhoff
Inria Saclay – Ile-de-France

Course Overview

1	Mon, 18.9.2017 Tue, 19.9.2017	first lecture groups defined via wiki everybody went (actively!) through the Getting Started part of github.com/numbbo/coco
2	Wed, 20.9.2017	lecture: "Benchmarking", final adjustments of groups everybody can run and postprocess the example experiment (~1h for final questions/help during the lecture)
3	Fri, 22.9.2017	today's lecture "Introduction to Continuous Optimization"
4	Fri, 29.9.2017	lecture "Gradient-Based Algorithms"
5	Fri, 6.10.2017	lecture "Stochastic Algorithms and DFO"
6	Fri, 13.10.2017	lecture "Discrete Optimization I: graphs, greedy algos, dyn. progr." deadline for submitting data sets
	Wed, 18.10.2017	deadline for paper submission
7	Fri, 20.10.2017	final lecture "Discrete Optimization II: dyn. progr., B&B, heuristics"
	Thu, 26.10.2017 / Fri, 27.10.2017	oral presentations (individual time slots)
	after 30.10.2017	vacation aka learning for the exams
	Fri, 10.11.2017	written exam

All deadlines:
23:59pm Paris time

Details on Continuous Optimization Lectures

Introduction to Continuous Optimization

- examples (from ML / black-box problems)
- typical difficulties in optimization

Mathematical Tools to Characterize Optima

- reminders about differentiability, gradient, Hessian matrix
 - unconstraint optimization
 - first and second order conditions
 - convexity
 - constraint optimization
-

Gradient-based Algorithms

- quasi-Newton method (BFGS)

DFO: trust-region method (Nelder-Mead)

Learning in Optimization / Stochastic Optimization

- CMA-ES (adaptive algorithms / Information Geometry)
- PhD thesis possible on this topic

*method strongly related to ML / new promising research area
interesting open questions*

CMA-ES in a Nutshell

The CMA-ES

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

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

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

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

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

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

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

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

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

Not covered on this slide: termination
encoding

Goal:

Understand the main principles
of this state-of-the-art algorithm.

Copyright Notice

- Last slide was taken from
<https://www.lri.fr/~hansen/copenhagen-cma-es.pdf>
(copyright by Nikolaus Hansen, one of the main inventors of the CMA-ES algorithms)
- In the following, I will borrow more slides from there and from
<http://researchers.lille.inria.fr/~brockhoff/optimizationSaclay/slides/20151106-continuousoptIV.pdf>
(by Anne Auger)
- In the following and the online material in particular, I refer to these pdfs as [Hansen, p. X] and [Auger, p. Y] respectively.

The CMA-ES

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

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

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

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

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

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

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

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

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

Not covered on this slide: termination
encoding

Goal:

Understand the main principles
of this state-of-the-art algorithm.

CMA-ES: Stochastic Search Template

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

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

While happy do:

- Sample distribution $P(x|\theta) \rightarrow x_1, \dots, x_\lambda \in \mathbb{R}^n$
- Evaluate x_1, \dots, x_λ on f
- Update parameters $\theta \leftarrow F_\theta(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

For CMA-ES and evolution strategies in general:

sample distributions = multivariate Gaussian distributions

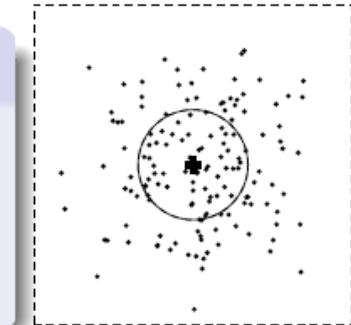
Sampling New Candidate Solutions (Offspring)

Evolution Strategies

New search points are sampled normally distributed

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

as perturbations of \mathbf{m} , where $x_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$



where

- the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ 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

here, all new points are sampled with the same parameters

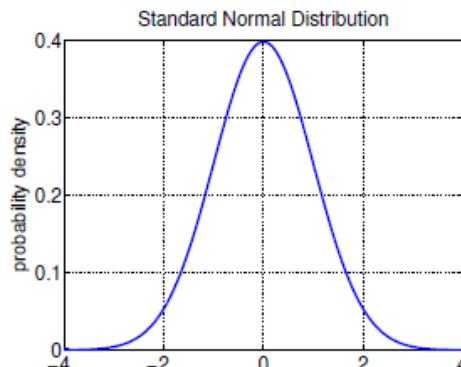
it remains to show how to adapt the parameters, but for now: normal distributions

from [Auger, p. 10]

Excursion: Normal Distributions

Normal Distribution

1-D case



probability density of the 1-D standard normal distribution $\mathcal{N}(0, 1)$
(expected (mean) value, variance) = (0, 1)

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

General case

- ▶ Normal distribution $\mathcal{N}(\mathbf{m}, \sigma^2)$
(expected value, variance) = (\mathbf{m}, σ^2)
density: $p_{\mathbf{m}, \sigma}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mathbf{m})^2}{2\sigma^2}\right)$
- ▶ A normal distribution is entirely determined by its mean value and variance
- ▶ The family of normal distributions is closed under linear transformations: if X is normally distributed then a linear transformation $aX + b$ is also normally distributed
- ▶ **Exercice:** Show that $\mathbf{m} + \sigma\mathcal{N}(0, 1) = \mathcal{N}(\mathbf{m}, \sigma^2)$

from [Auger, p. 11]

Excursion: Normal Distributions

Normal Distribution

General case

A random variable following a 1-D normal distribution is determined by its mean value m and variance σ^2 .

In the n -dimensional case it is determined by its mean vector and covariance matrix

Covariance Matrix

If the entries in a vector $X = (X_1, \dots, X_n)^T$ are random variables, each with finite variance, then the covariance matrix Σ is the matrix whose (i, j) entries are the covariance of (X_i, X_j)

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$$

where $\mu_i = E(X_i)$. Considering the expectation of a matrix as the expectation of each entry, we have

$$\Sigma = E[(X - \mu)(X - \mu)^T]$$

Σ is symmetric, positive definite
from [Auger, p. 12]

Excursion: Normal Distributions

The Multi-Variate (n -Dimensional) Normal Distribution

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

density: $p_{\mathcal{N}(\mathbf{m}, \mathbf{C})}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})\right),$

from [Auger, p. 13]

Excursion: Normal Distributions

The Multi-Variate (n -Dimensional) Normal Distribution

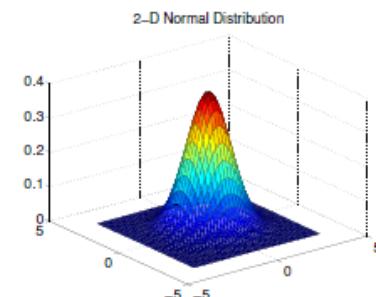
Any multi-variate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is uniquely determined by its mean value $\mathbf{m} \in \mathbb{R}^n$ and its symmetric positive definite $n \times n$ covariance matrix \mathbf{C} .

$$\text{density: } p_{\mathcal{N}(\mathbf{m}, \mathbf{C})}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})\right),$$

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

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

$$\mathcal{N}(\mathbf{m}, \mathbf{C}) = \mathbf{m} + \mathcal{N}(0, \mathbf{C})$$



from [Auger, p. 13]

Excursion: Normal Distributions

The Multi-Variate (n -Dimensional) Normal Distribution

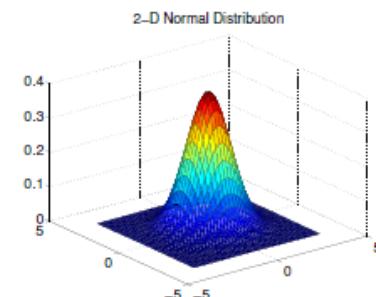
Any multi-variate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is uniquely determined by its mean value $\mathbf{m} \in \mathbb{R}^n$ and its symmetric positive definite $n \times n$ covariance matrix \mathbf{C} .

$$\text{density: } p_{\mathcal{N}(\mathbf{m}, \mathbf{C})}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})\right),$$

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

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

$$\mathcal{N}(\mathbf{m}, \mathbf{C}) = \mathbf{m} + \mathcal{N}(0, \mathbf{C})$$



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

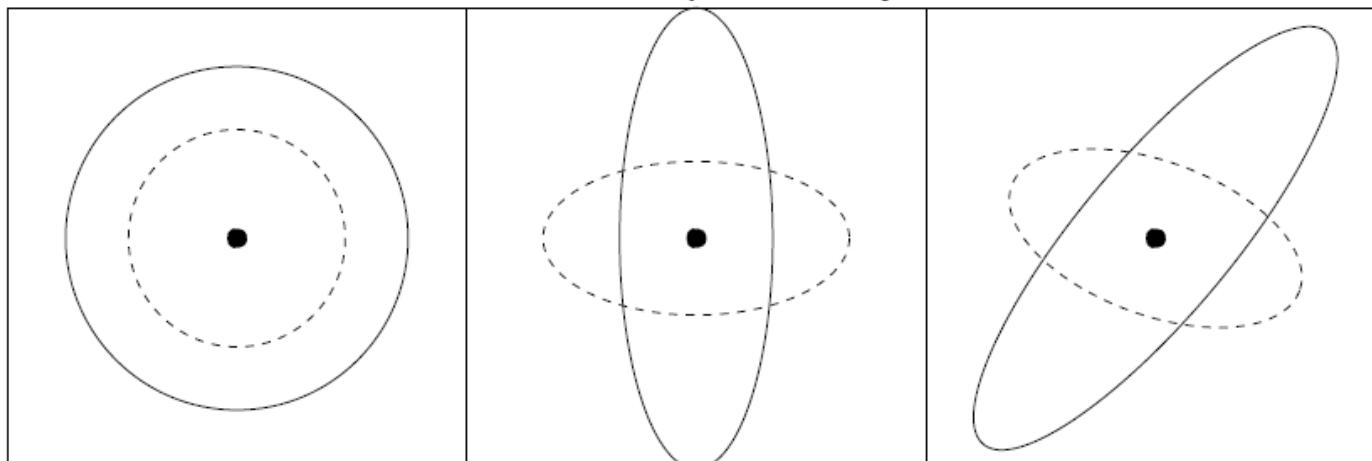
- ▶ determines the shape
- ▶ **geometrical interpretation:** any covariance matrix can be uniquely identified with the iso-density ellipsoid
$$\{\mathbf{x} \in \mathbb{R}^n \mid (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}) = 1\}$$

from [Auger, p. 13]

Covariance Matrix: Lines of Equal Density

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

Lines of Equal Density



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

one degree of freedom σ

components are
independent standard
normally distributed

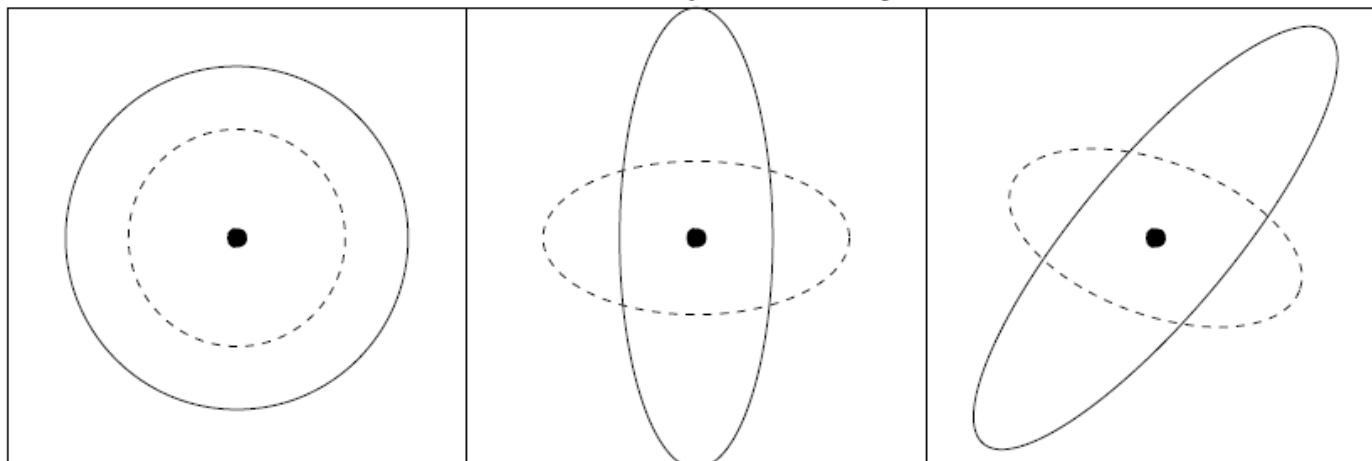
where \mathbf{I} is the identity matrix (isotropic case) and \mathbf{D} is a diagonal matrix (reasonable for separable problems) and $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T)$ holds for all \mathbf{A} .

from [Auger, p. 14]

Covariance Matrix: Lines of Equal Density

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

Lines of Equal Density



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

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

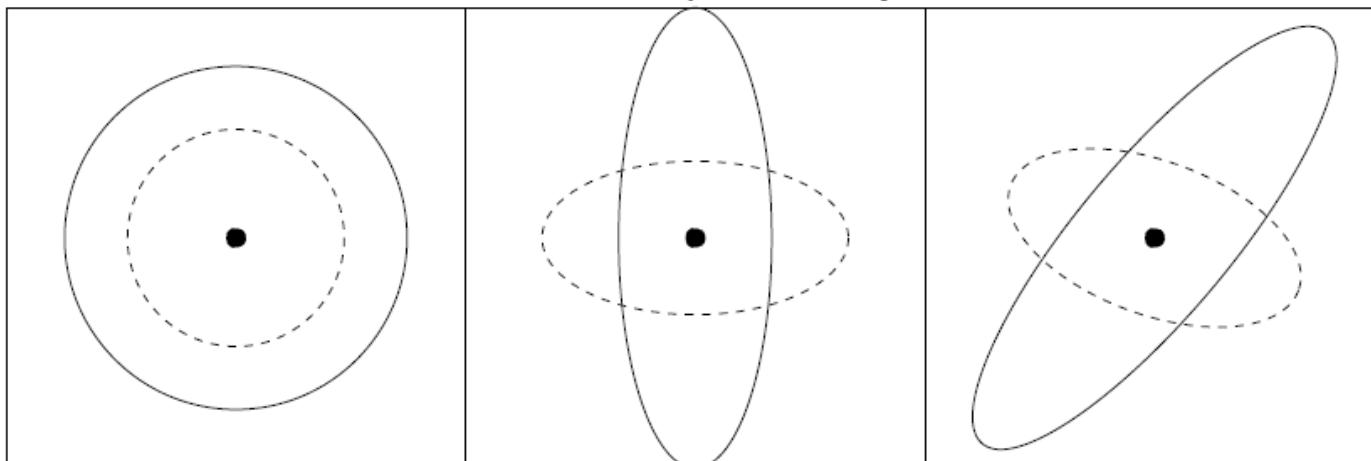
where \mathbf{I} is the identity matrix (isotropic case) and \mathbf{D} is a diagonal matrix (reasonable for separable problems) and $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T)$ holds for all \mathbf{A} .

from [Auger, p. 14]

Covariance Matrix: Lines of Equal Density

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

Lines of Equal Density



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

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

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

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

from [Auger, p. 14]

Adaptation of Sample Distribution Parameters

Adaptation: What do we want to achieve?

New search points are sampled normally distributed

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

where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$

- ▶ the **mean** vector should represent the favorite solution
 - ▶ the **step-size** controls the step-length and thus convergence rate
 - should allow to reach fastest convergence rate possible
 - ▶ the **covariance matrix** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid
 - adaptation should allow to learn the “topography” of the problem
 - particulary important for ill-conditionned problems
 - $\mathbf{C} \propto \mathbf{H}^{-1}$ on convex quadratic functions
- from [Auger, p. 16]

Adaptation of the Mean

Plus and Comma Selection

Evolution Strategies

Terminology

μ : # of parents, λ : # of offspring

Plus (elitist) and comma (non-elitist) selection

$(\mu + \lambda)$ -ES: selection in $\{\text{parents}\} \cup \{\text{offspring}\}$

(μ, λ) -ES: selection in $\{\text{offspring}\}$

$(1 + 1)$ -ES

Sample one offspring from parent $\textcolor{green}{m}$

$$\boldsymbol{x} = \textcolor{green}{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C})$$

If \boldsymbol{x} better than $\textcolor{green}{m}$ select

$$\textcolor{green}{m} \leftarrow \boldsymbol{x}$$

Non-Elitism and Weighted Recombination

The $(\mu/\mu, \lambda)$ -ES

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

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

where

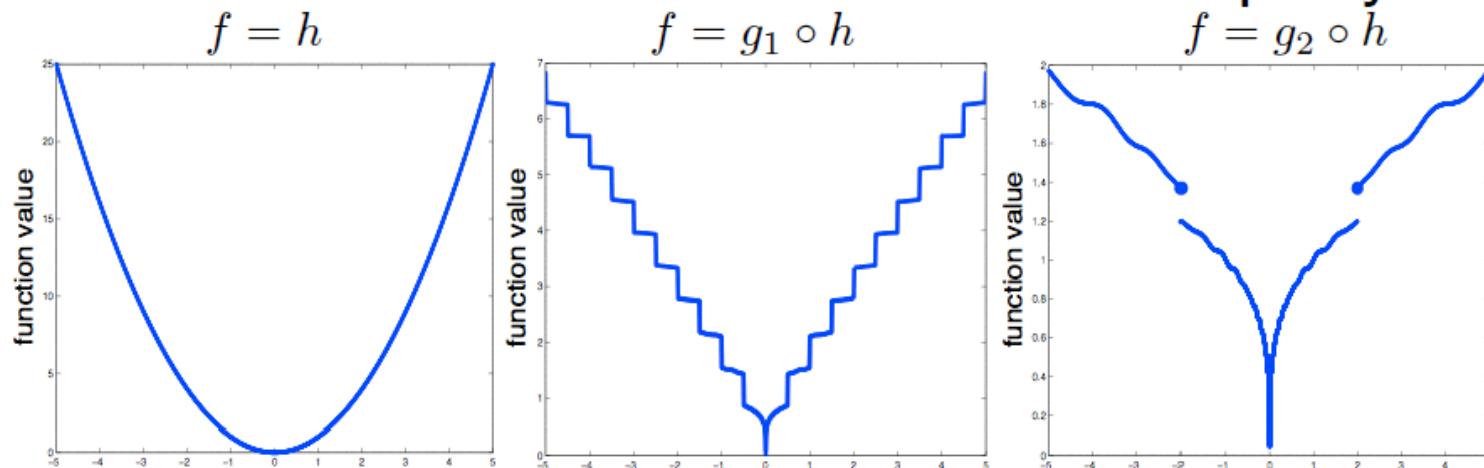
$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$

The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

from [Hansen, p. 34]
26 / 81

Invariance Against Order-Preserving f -Transformations

Invariance: Function-Value Free Property



Three functions belonging to the same equivalence class

A *function-value free search algorithm* is invariant under the transformation with any **order preserving** (strictly increasing) g .

Invariances make

- observations meaningful as a rigorous notion of generalization
- algorithms predictable and/or "robust"

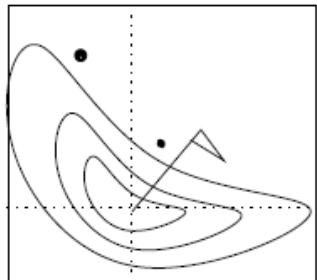
from [Hansen, p. 37]

Invariance Against Translations in Search Space

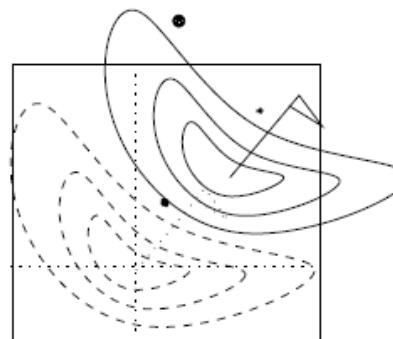
Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



$$f(\mathbf{x}) \leftrightarrow f(\mathbf{x} - \mathbf{a})$$



Identical behavior on f and f_a

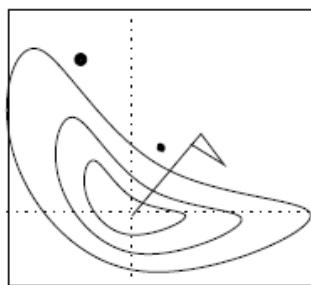
$$\begin{aligned} f : \quad \mathbf{x} &\mapsto f(\mathbf{x}), & \mathbf{x}^{(t=0)} &= \mathbf{x}_0 \\ f_a : \quad \mathbf{x} &\mapsto f(\mathbf{x} - \mathbf{a}), & \mathbf{x}^{(t=0)} &= \mathbf{x}_0 + \mathbf{a} \end{aligned}$$

No difference can be observed w.r.t. the argument of f

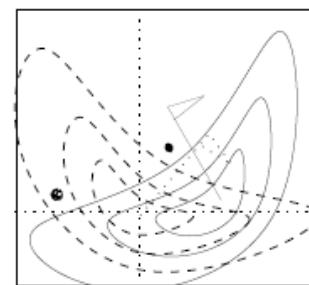
Invariance Against Search Space Rotations

Rotational Invariance in Search Space

- invariance to orthogonal (rigid) transformations \mathbf{R} , where $\mathbf{RR}^T = \mathbf{I}$
e.g. true for simple evolution strategies
recombination operators might jeopardize rotational invariance



$$f(\mathbf{x}) \leftrightarrow f(\mathbf{Rx})$$



Identical behavior on f and $f_{\mathbf{R}}$

$$\begin{aligned} f : \quad \mathbf{x} &\mapsto f(\mathbf{x}), & \mathbf{x}^{(t=0)} &= \mathbf{x}_0 \\ f_{\mathbf{R}} : \quad \mathbf{x} &\mapsto f(\mathbf{Rx}), & \mathbf{x}^{(t=0)} &= \mathbf{R}^{-1}(\mathbf{x}_0) \end{aligned}$$

45

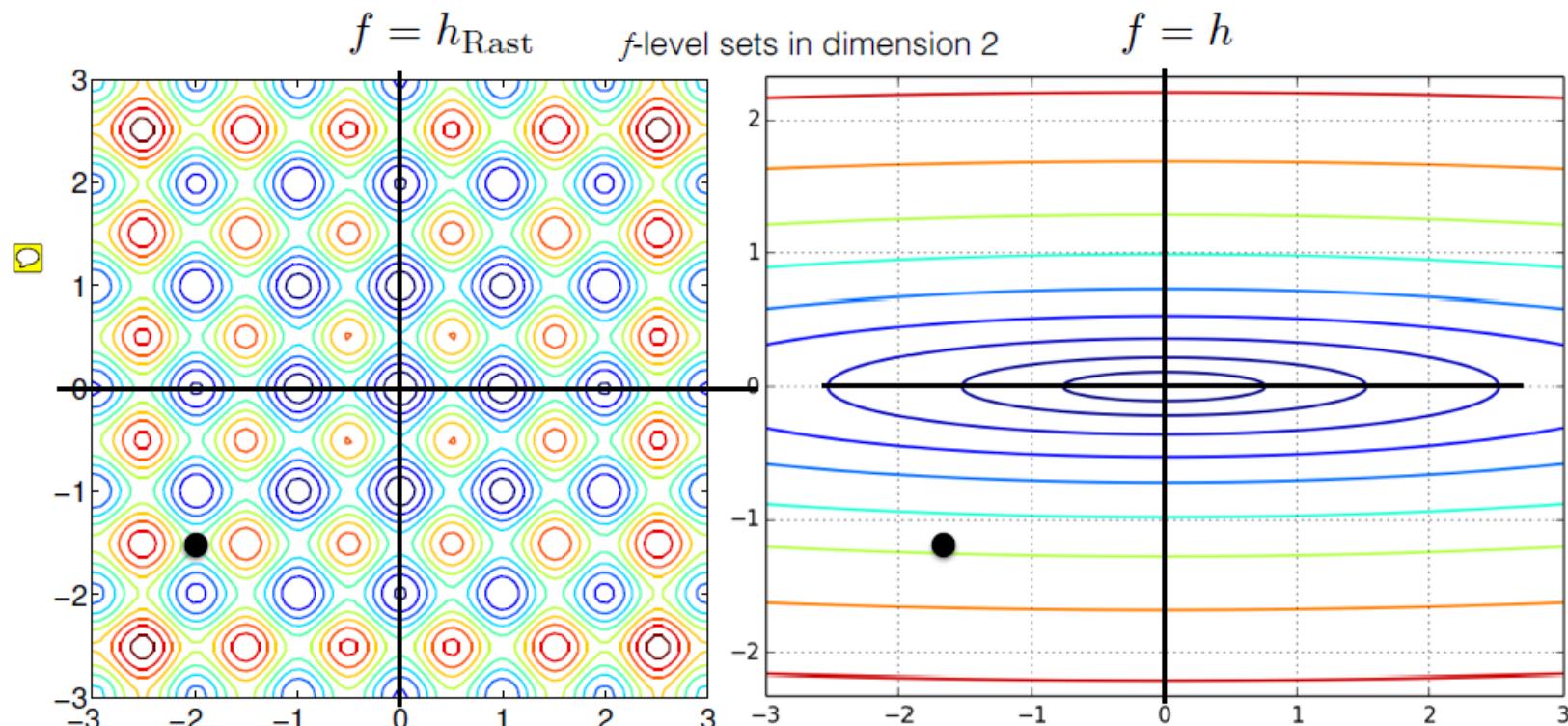
No difference can be observed w.r.t. the argument of f

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

⁵ Hansen 2000. Invariance, Self-Adaptation and Correlated Mutations in Evolution Strategies. *Parallel Problem Solving from Nature PPSN VI*

Invariance Against Rigid Search Space Transformations

Invariance Under Rigid Search Space Transformations

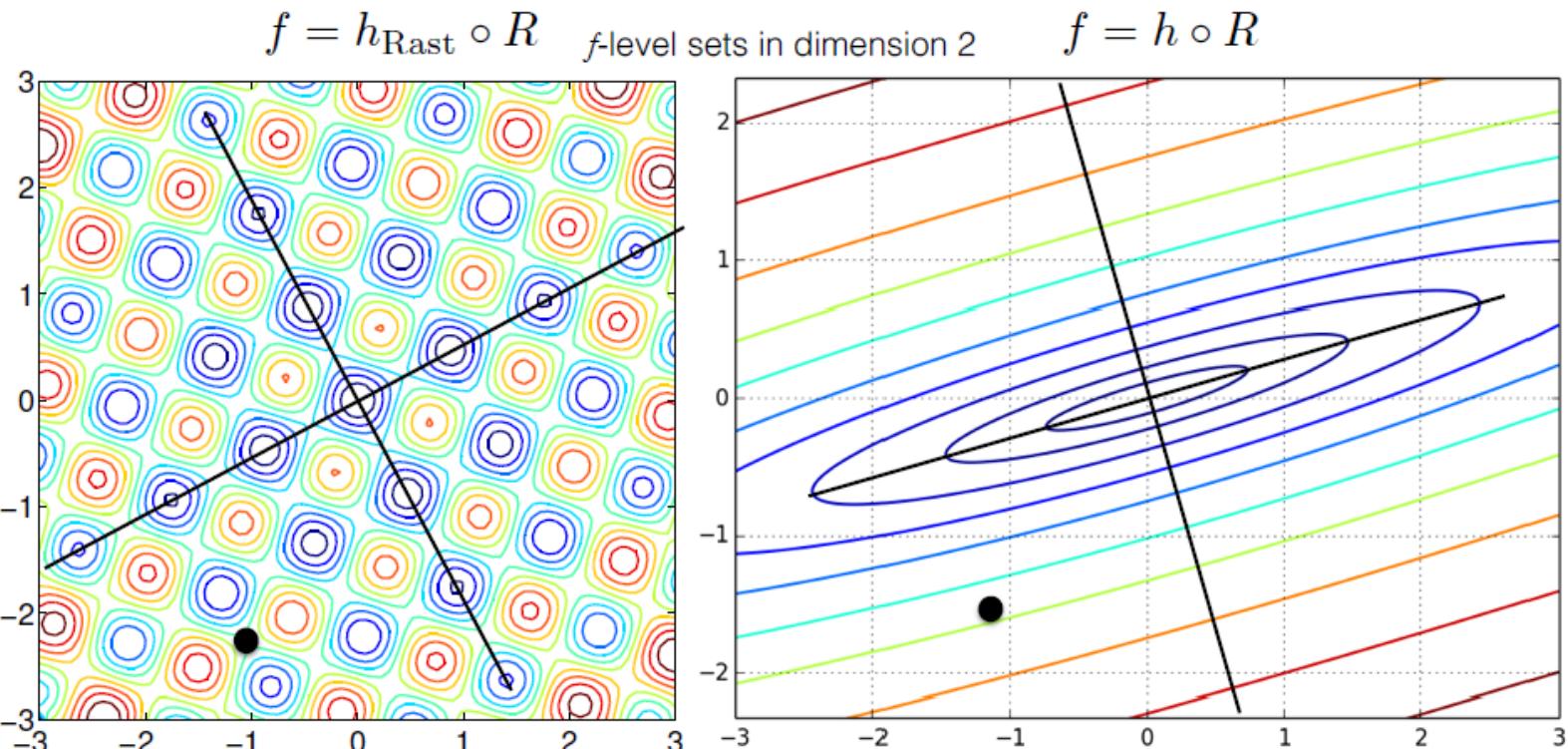


for example, invariance under search space rotation
(separable \Leftrightarrow non-separable)

from [Hansen, p. 40]

Invariance Against Rigid Search Space Transformations

Invariance Under Rigid Search Space Transformations



for example, invariance under search space rotation
(separable \Leftrightarrow non-separable)

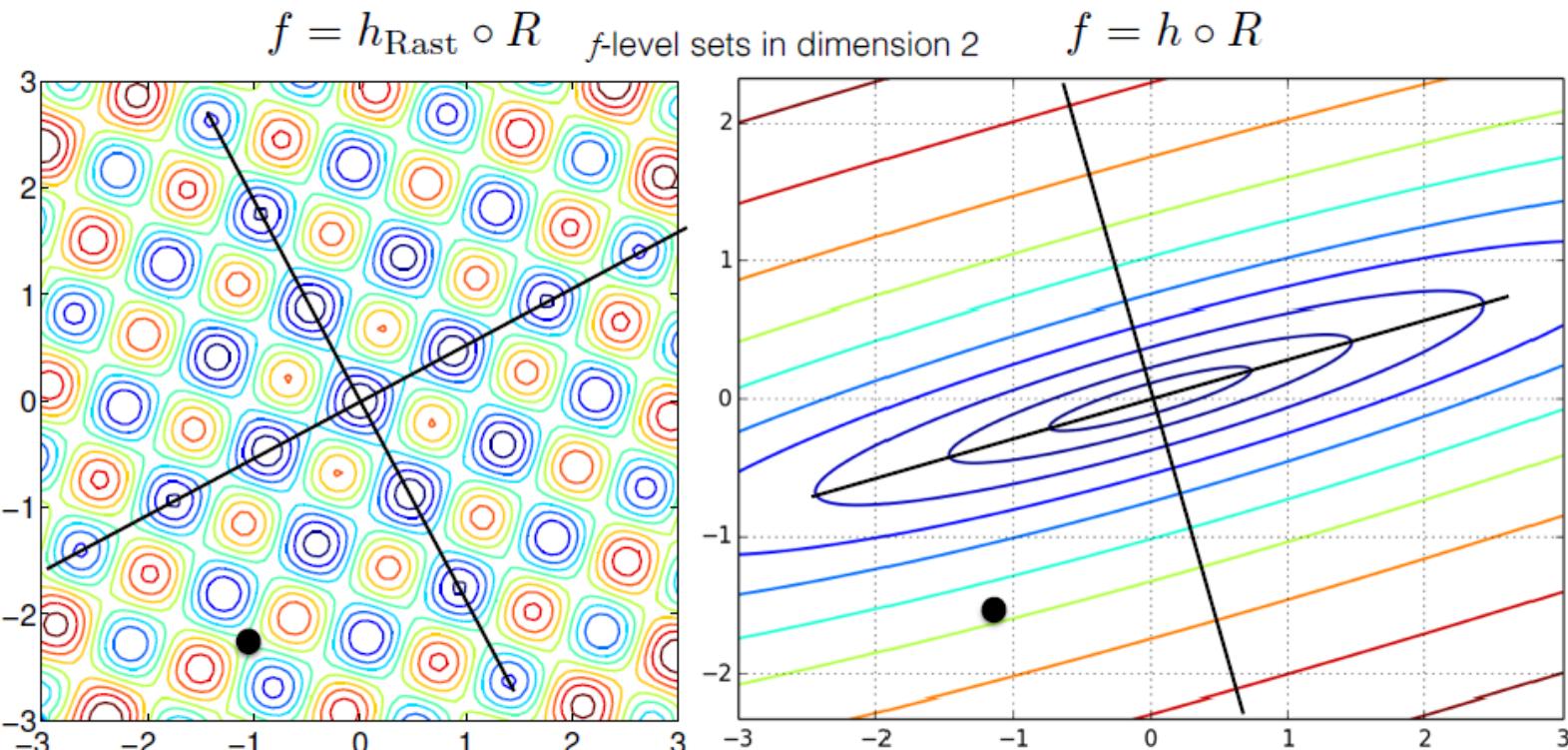
from [Hansen, p. 41]

Invariance Against Rigid Search Space Transformations

Evolution Strategies (ES)

Invariance

Invariance Under Rigid Search Space Transformations



for example, invariance under
(separable \Leftrightarrow non-separable)

mainly Nelder-Mead and CMA-ES
have this property

Invariances: Summary

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

- Empirical performance results

- ▶ from benchmark functions
 - ▶ from solved real world problems

are only useful if they do **generalize** to other problems

- **Invariance** is a strong **non-empirical** statement about generalization

generalizing (identical) performance from a single function to a whole class of functions

consequently, invariance is important for the evaluation of search algorithms

Step-Size Adaptation

Recap CMA-ES: What We Have So Far

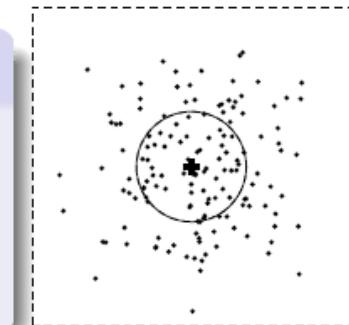
Evolution Strategies

Recalling

New search points are sampled normally distributed

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

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$



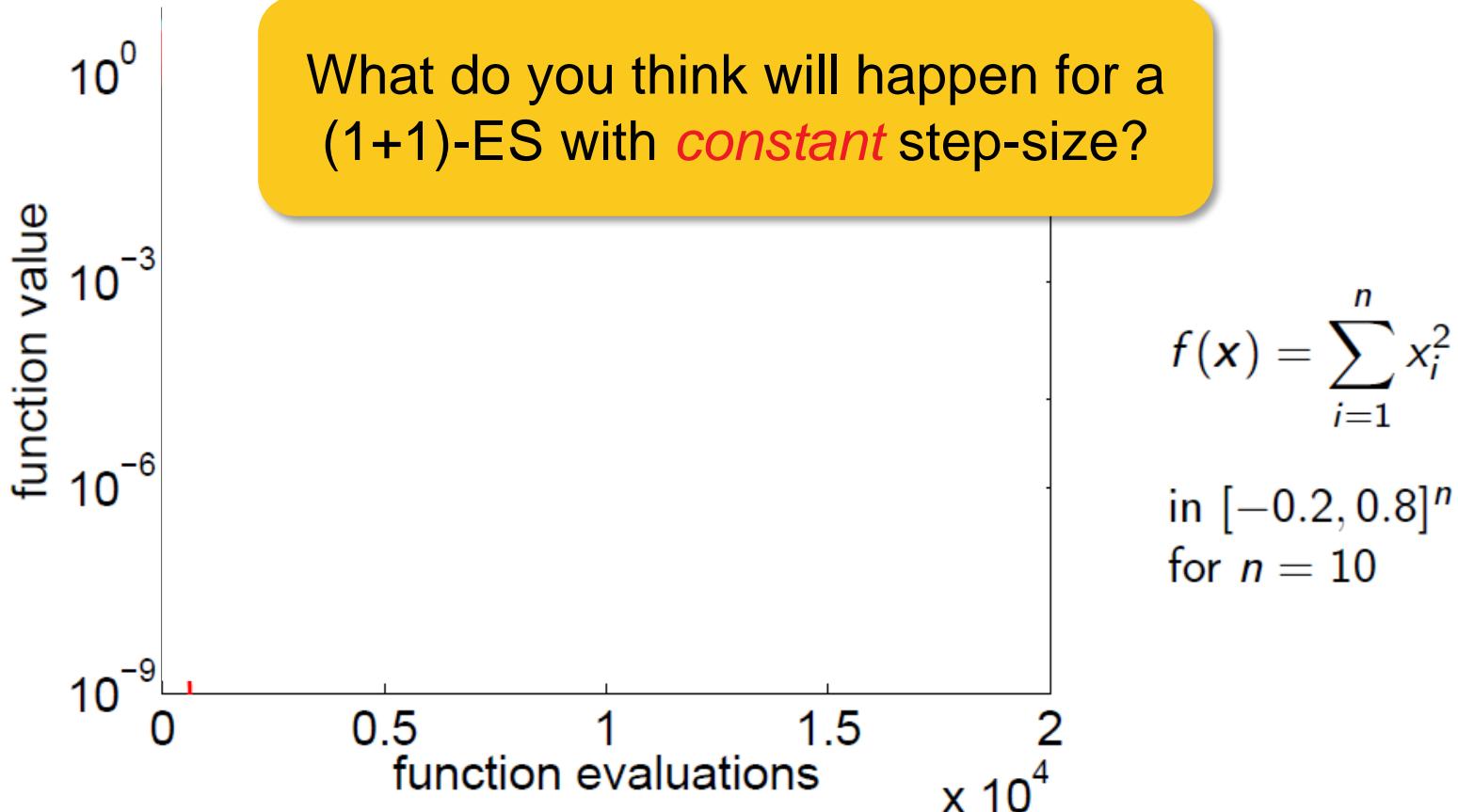
where

- the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution and $\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$
- 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 remaining question is how to update σ and \mathbf{C} .

Why At All Step-Size Adaptation?

Why Step-Size Control?

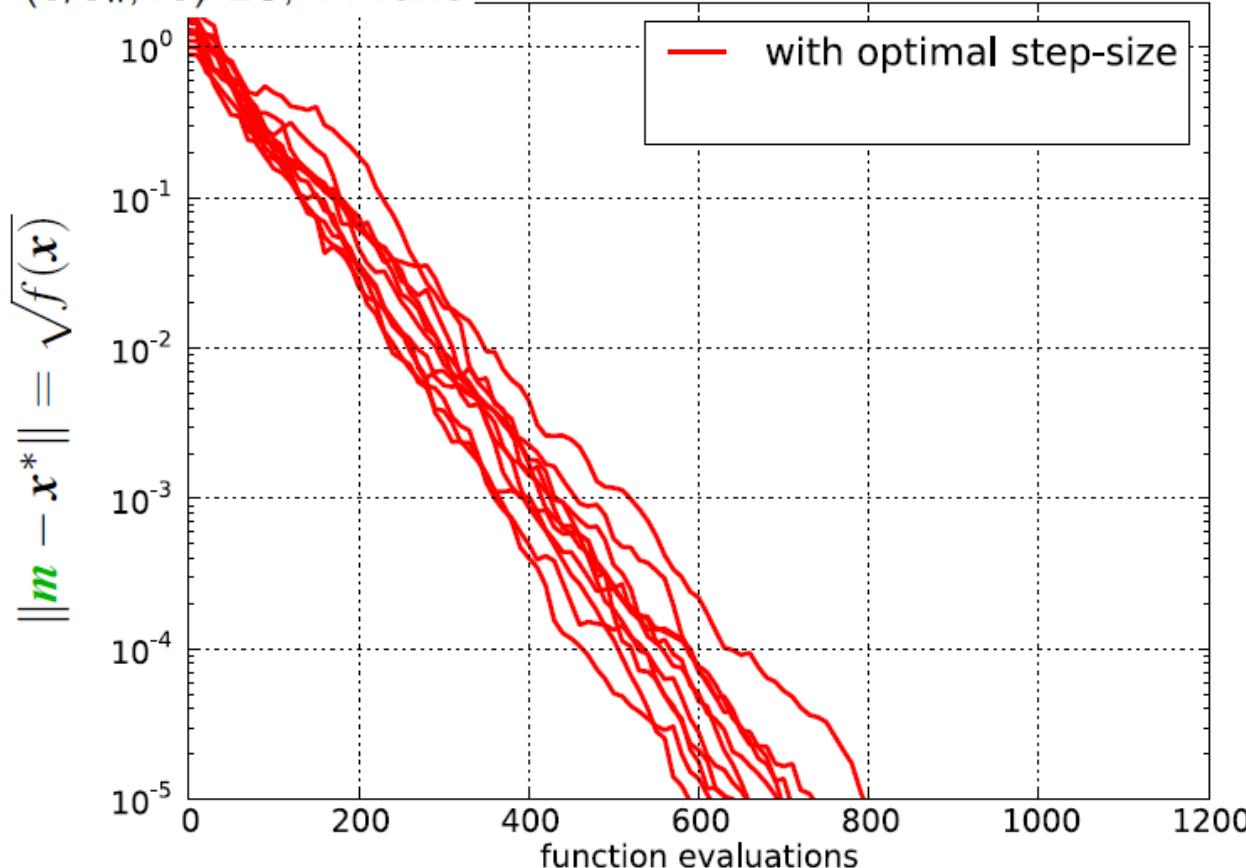


from [Auger, p. 22]

Optimal Step-Size

Why Step-Size Control?

(5/5_w, 10)-ES, 11 runs



with optimal step-size σ

$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

from [Hansen, p. 47]

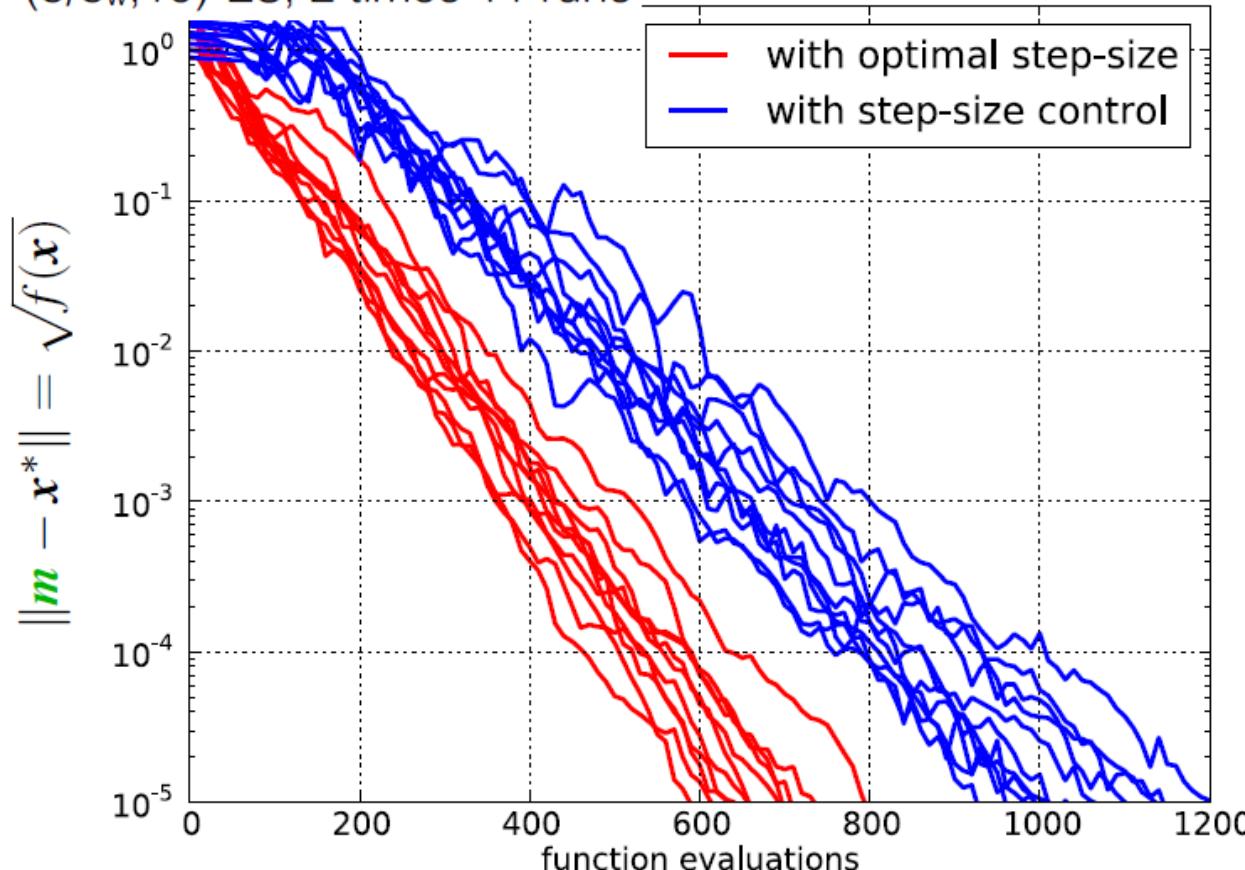
Optimal Step-Size vs. Step-Size Control

Step-Size Control

Why Step-Size Control

Why Step-Size Control?

(5/5_w, 10)-ES, 2 times 11 runs



$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

with **optimal** versus **adaptive** step-size σ with too small initial σ

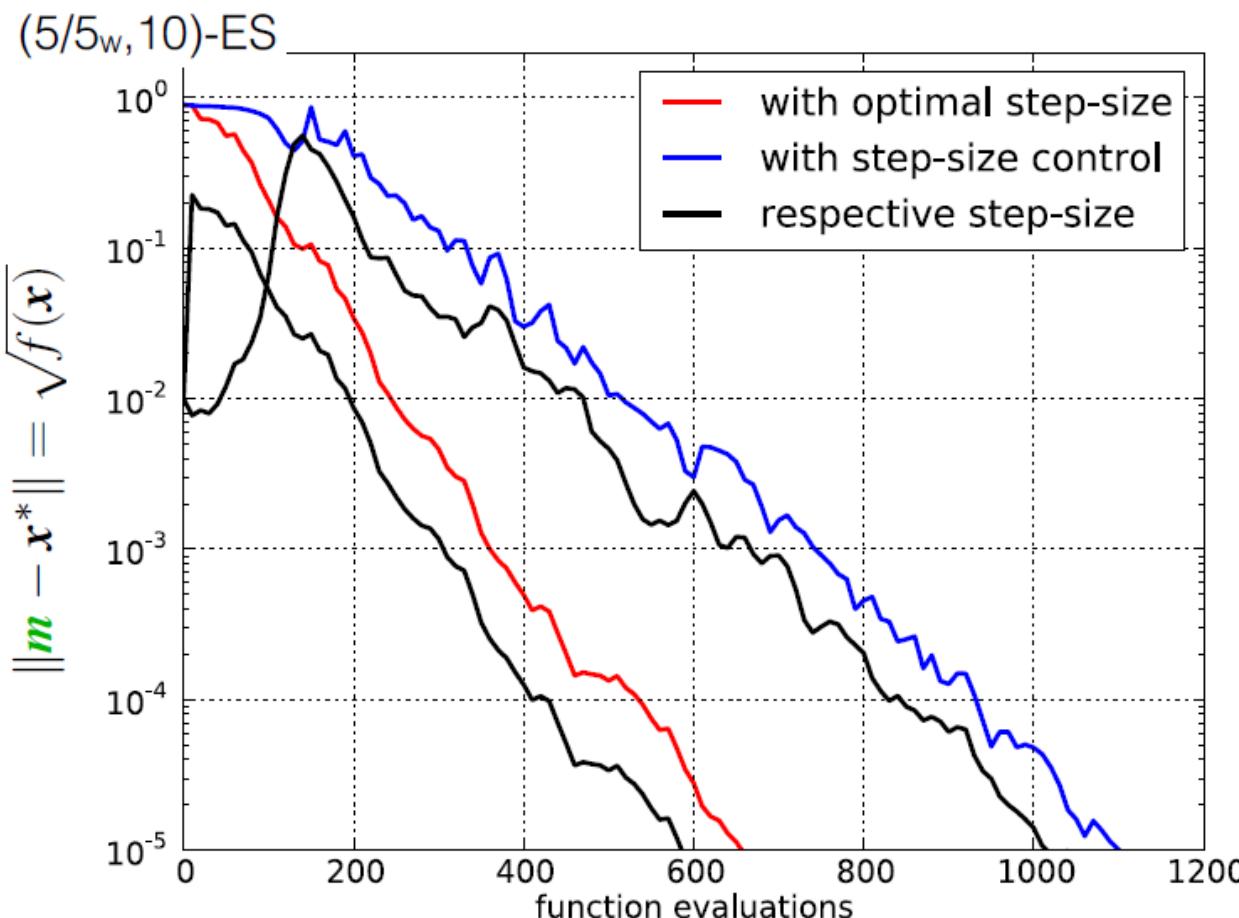
from [Hansen, p. 48]

Optimal Step-Size vs. Step-Size Control

Step-Size Control

Why Step-Size Control

Why Step-Size Control?



$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$ and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

comparing number of f -evals to reach $\|\mathbf{m}\| = 10^{-5}$: $\frac{1100 - 100}{650} \approx 1.5$

◀ □ ▶ ⏪ ⏩ ⏴ from [Hansen, p. 49]

Adapting the Step-Size

Question:

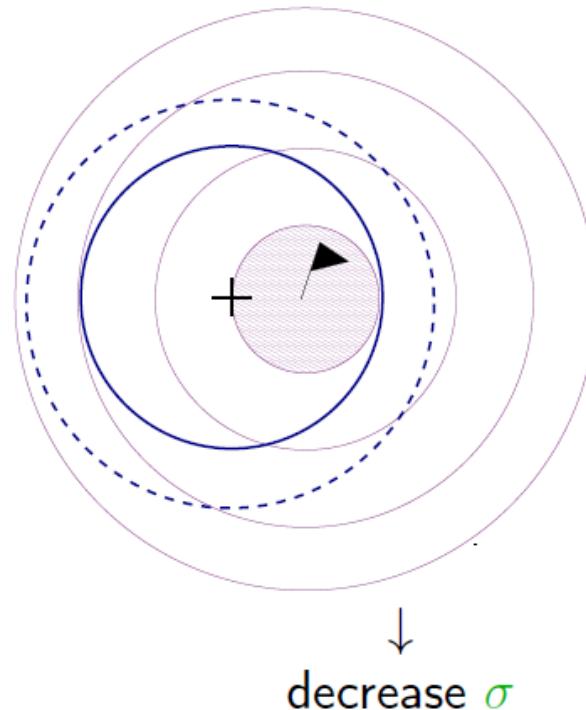
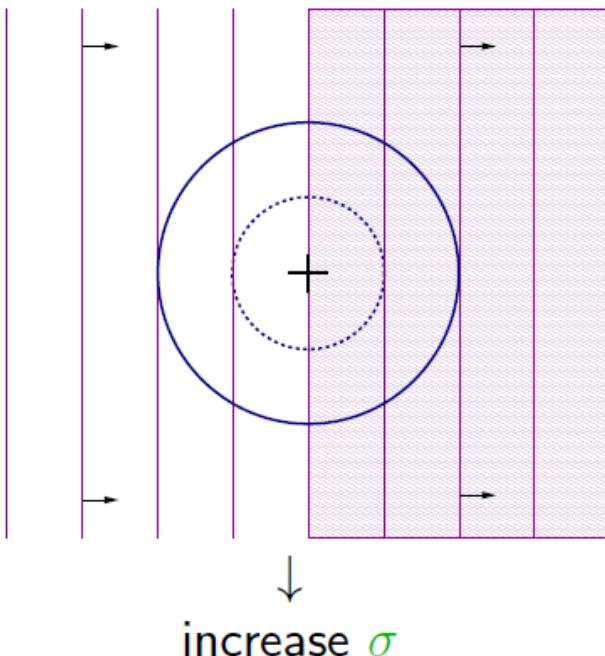
How to actually adapt the step-size during the optimization?

Most common:

- 1/5 success rule
- Cumulative Step-Size Adaptation (CSA, as in standard CMA-ES)
- others possible (Two-Point Adaptation, self-adaptive step-size, ...)

One-Fifth Success Rule

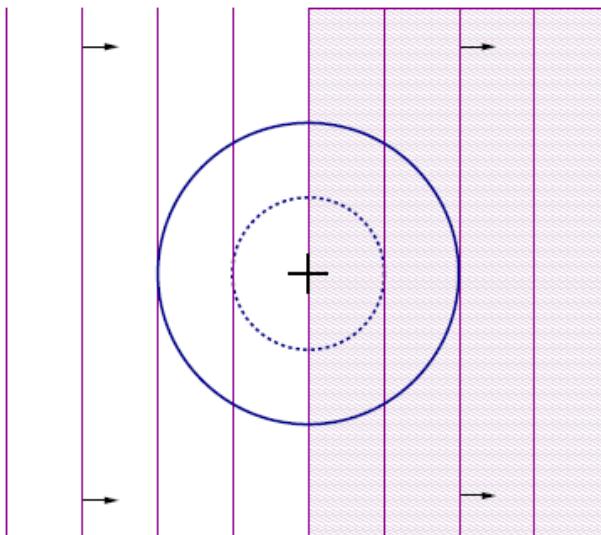
One-fifth success rule



from [Auger, p. 32]

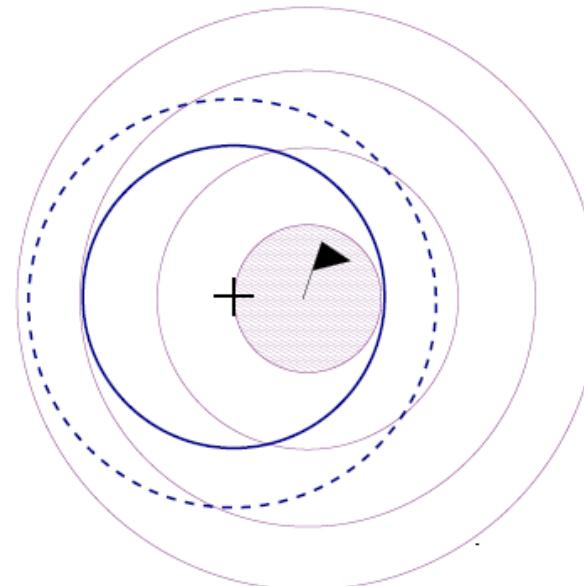
One-Fifth Success Rule

One-fifth success rule



Probability of success (p_s)

$1/2$



Probability of success (p_s)

“too small”

from [Auger, p. 33]

One-Fifth Success Rule

One-fifth success rule

p_s : # of successful offspring / # offspring (per generation)

$$\sigma \leftarrow \sigma \times \exp \left(\frac{1}{3} \times \frac{p_s - p_{\text{target}}}{1 - p_{\text{target}}} \right)$$

Increase σ if $p_s > p_{\text{target}}$
Decrease σ if $p_s < p_{\text{target}}$

(1 + 1)-ES

$$p_{\text{target}} = 1/5$$

IF offspring better parent

$$p_s = 1, \sigma \leftarrow \sigma \times \exp(1/3)$$

ELSE

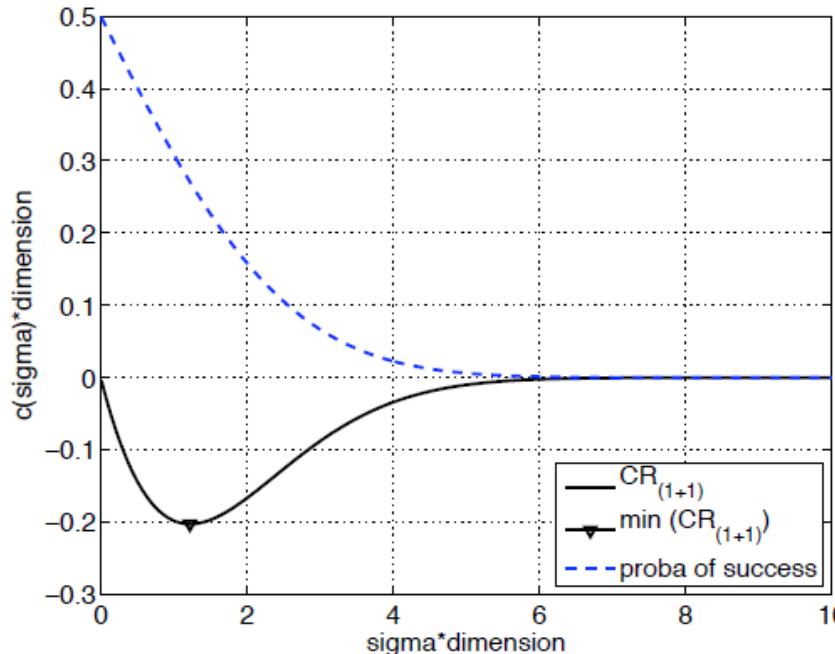
$$p_s = 0, \sigma \leftarrow \sigma / \exp(1/3)^{1/4}$$

from [Auger, p. 34]

One-Fifth Success Rule

Why 1/5?

Asymptotic convergence rate and probability of success of scale-invariant step-size (1+1)-ES



sphere - asymptotic results, i.e. $n = \infty$ (see slides before)

1/5 trade-off of optimal probability of success on the sphere and corridor

from [Auger, p. 35]

Cumulative Step-Size Adaptation (CSA)

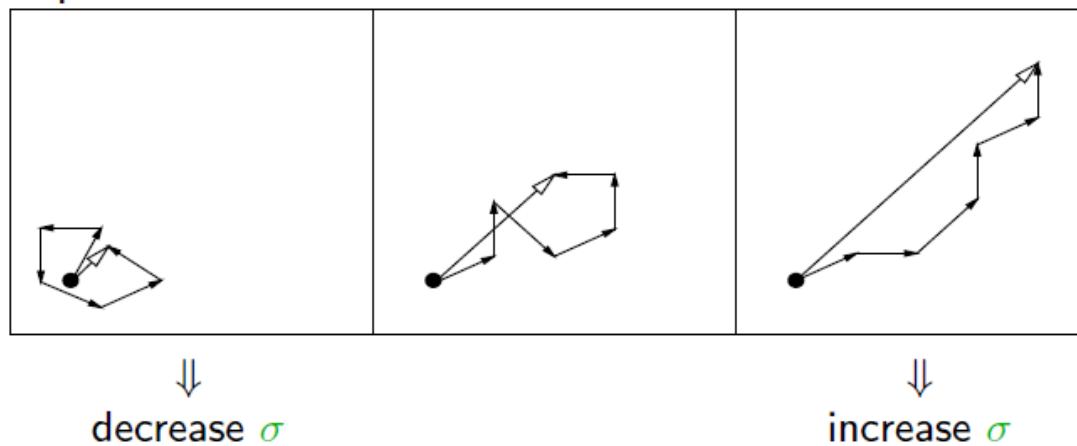
Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \end{aligned}$$

Measure the length of the *evolution path*

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



from [Auger, p. 36]

Cumulative Step-Size Adaptation (CSA)

Path Length Control (CSA)

The Equations

Initialize $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, evolution path $\mathbf{p}_\sigma = \mathbf{0}$,
set $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

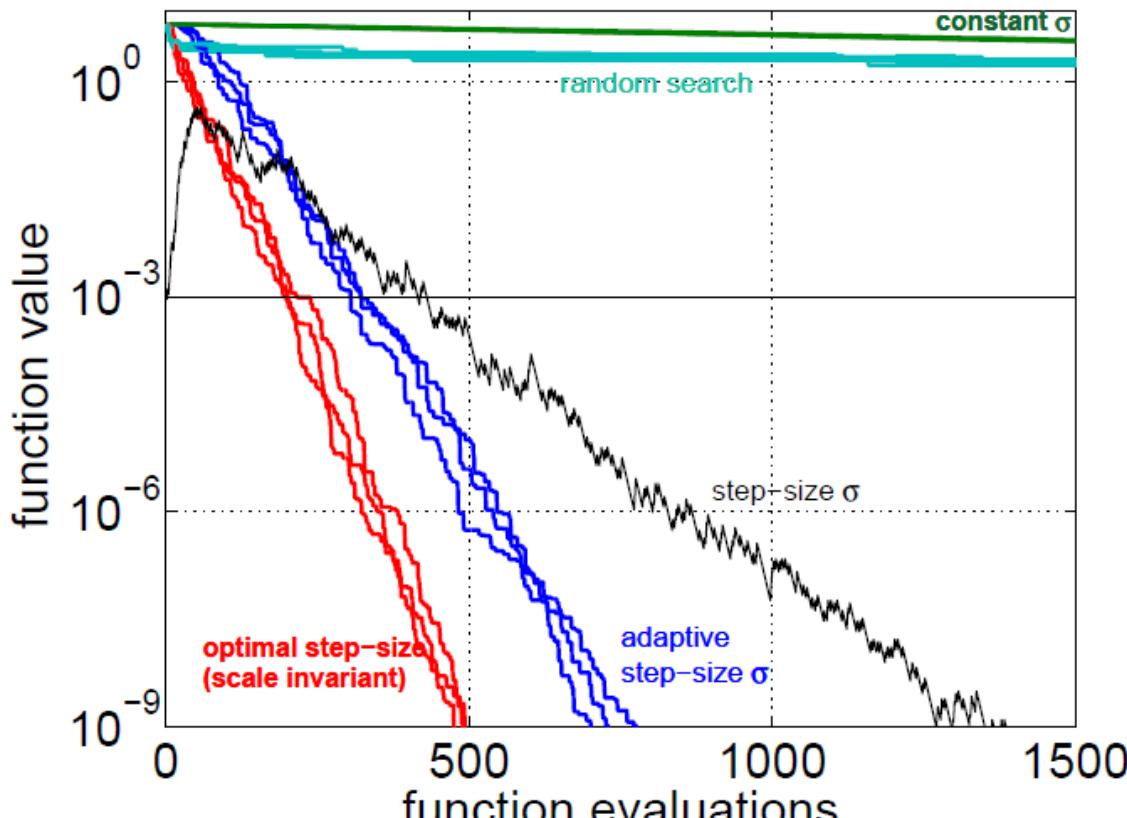
$$\begin{aligned}\mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} && \text{update mean} \\ \mathbf{p}_\sigma &\leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} \mathbf{y}_w \\ \sigma &\leftarrow \sigma \times \underbrace{\exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E} \|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)}_{>1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}} && \text{update step-size}\end{aligned}$$

from [Auger, p. 37]

Cumulative Step-Size Adaptation (CSA)

Step-size adaptation

What is achieved



$$f(x) = \sum_{i=1}^n x_i^2$$

in $[-0.2, 0.8]^n$
for $n = 10$

Linear convergence

from [Auger, p. 38]

Covariance Matrix Adaptation

Recap CMA-ES: What We Have So Far

Evolution Strategies

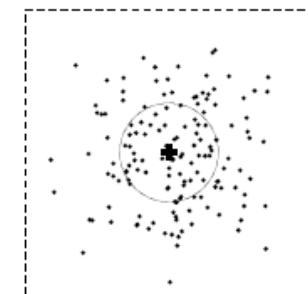
Recalling

New search points are sampled normally distributed

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

as perturbations of \mathbf{m} ,

where $x_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$,
 $\mathbf{C} \in \mathbb{R}^{n \times n}$



where

- ▶ the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ 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 remaining question is how to update \mathbf{C} .

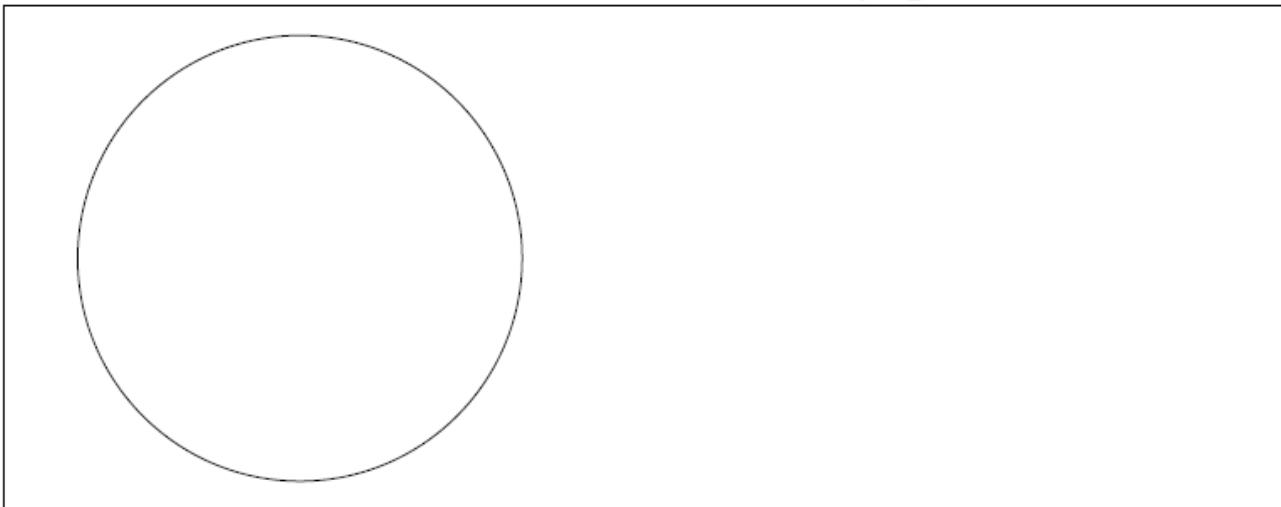
from [Auger, p. 40]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



initial distribution, $\mathbf{C} = \mathbf{I}$

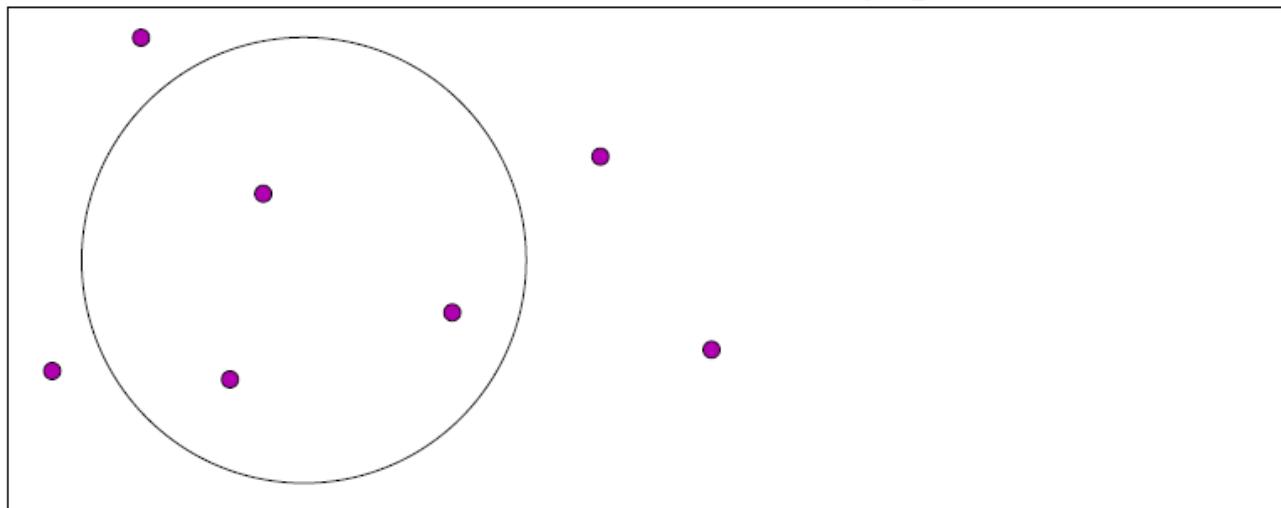
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



initial distribution, $\mathbf{C} = \mathbf{I}$

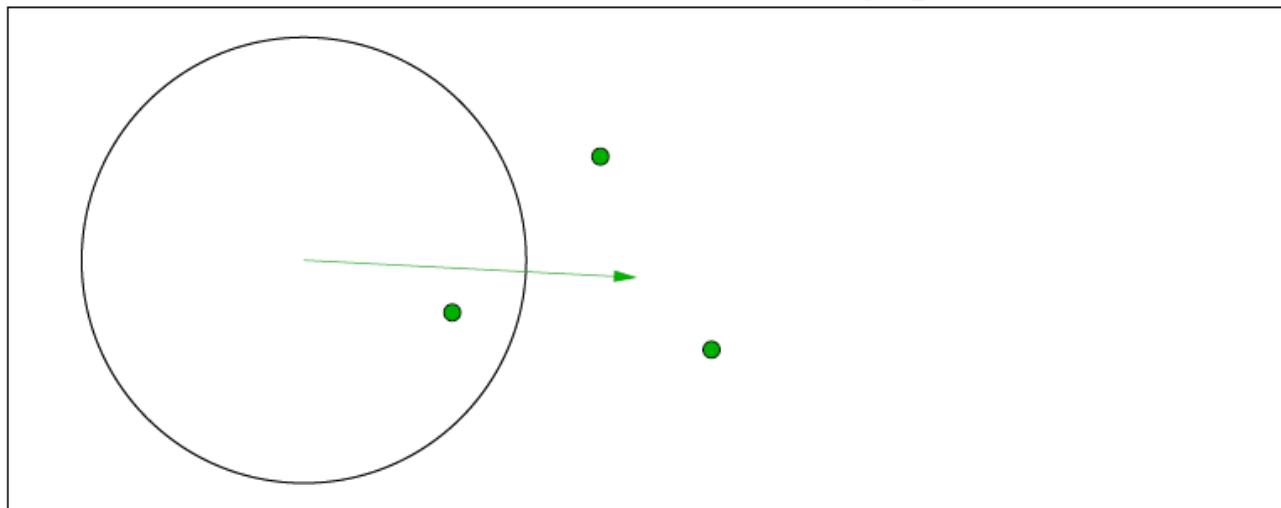
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



\mathbf{y}_w , movement of the population mean $\textcolor{blue}{m}$ (disregarding σ)

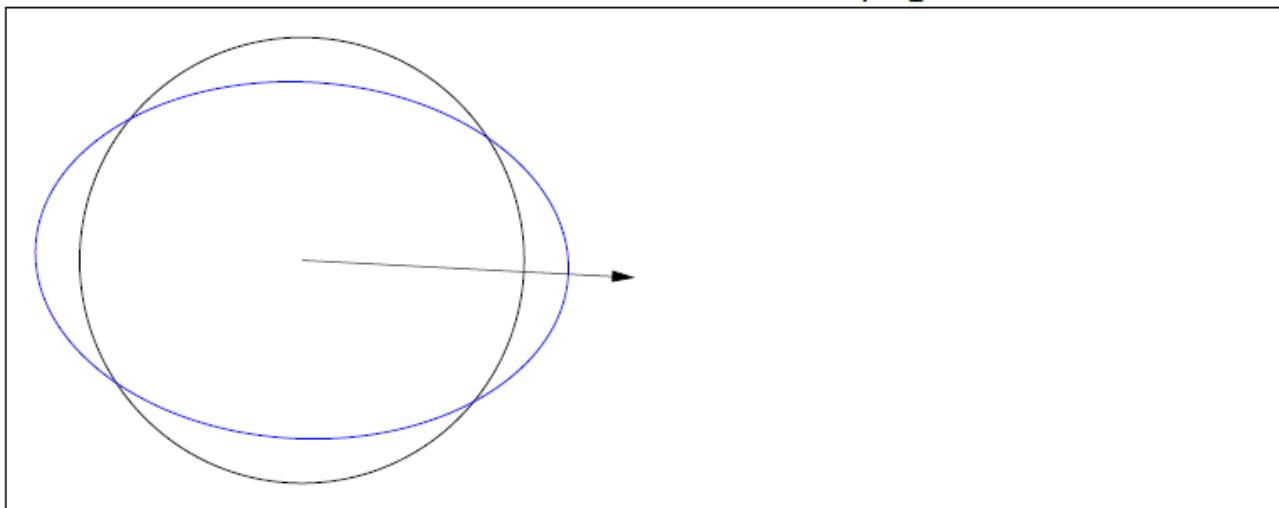
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

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

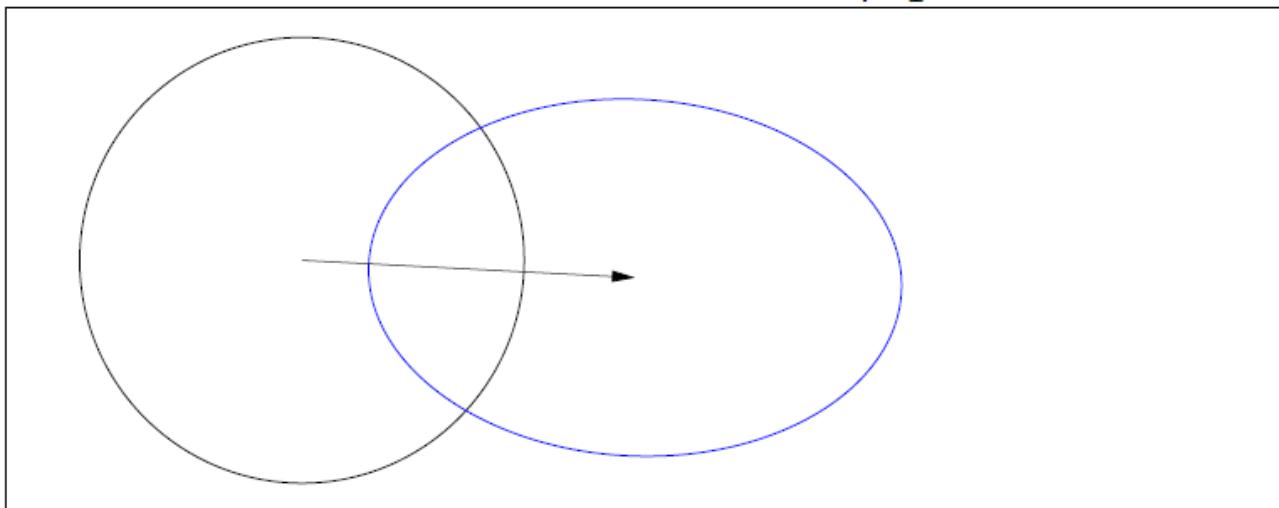
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



new distribution (disregarding σ)

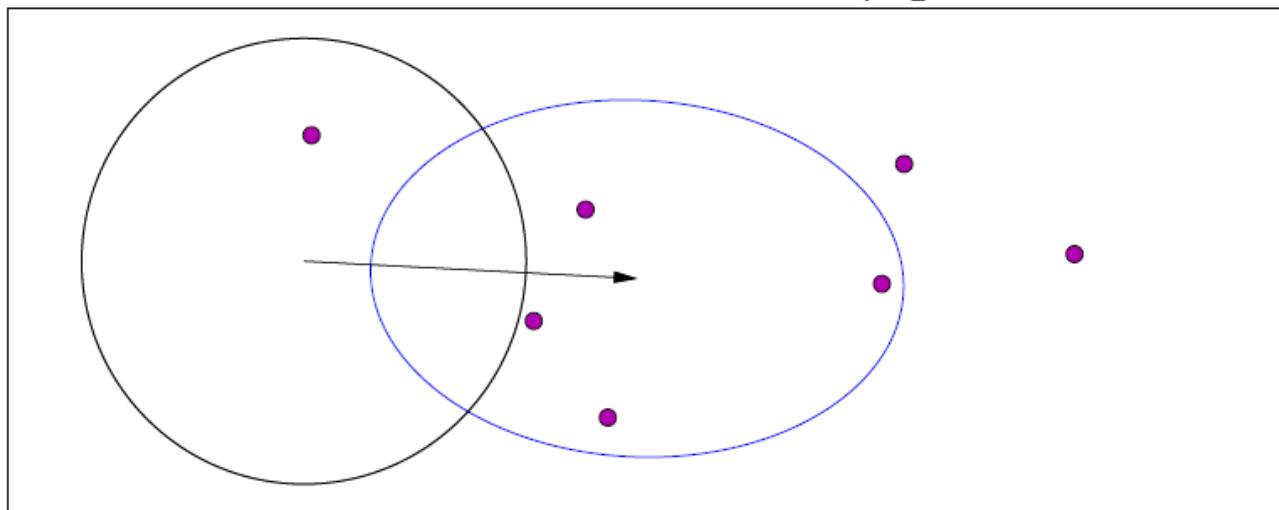
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



new distribution (disregarding σ)

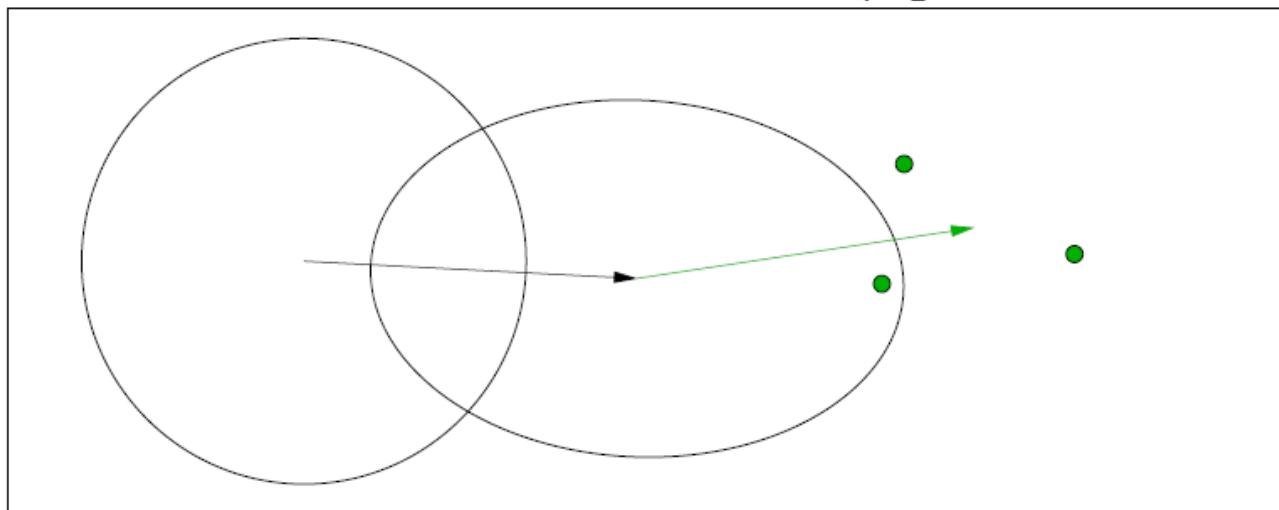
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



movement of the population mean \mathbf{m}

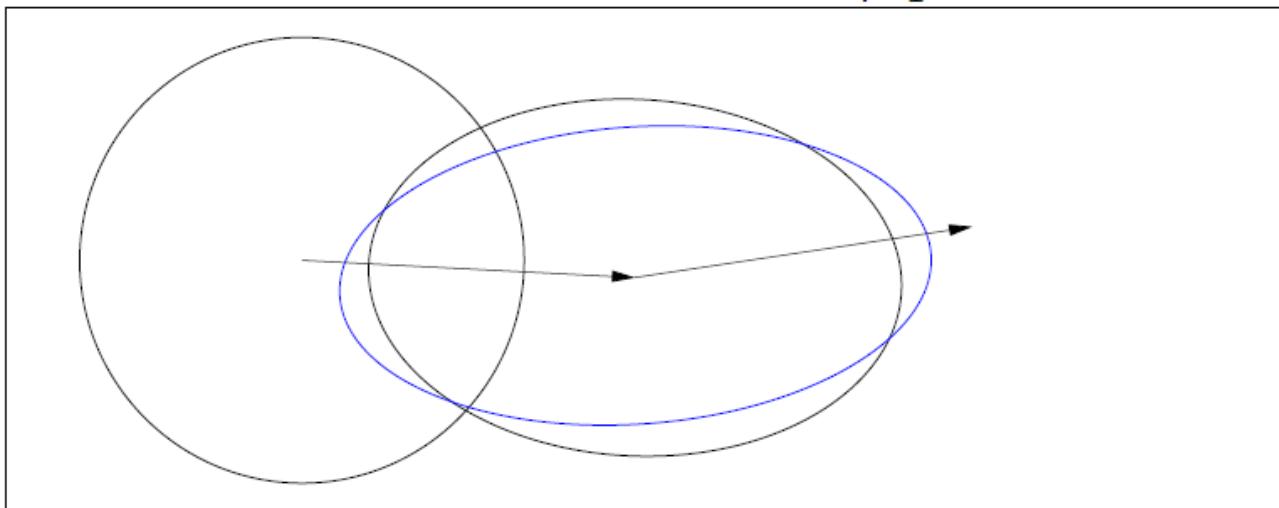
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

$$\textcolor{blue}{C} \leftarrow 0.8 \times \textcolor{blue}{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$$

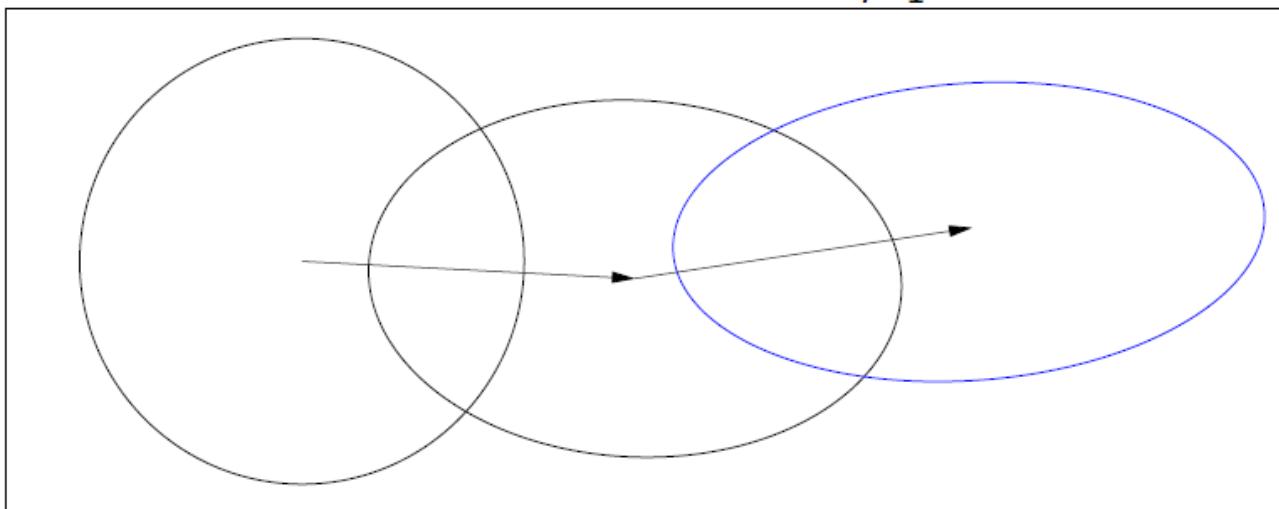
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



new distribution,

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

the ruling principle: the adaptation **increases the likelihood of successful steps**, \mathbf{y}_w , to appear again

from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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

While not terminate

$$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i, \quad \mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C}),$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda}$$

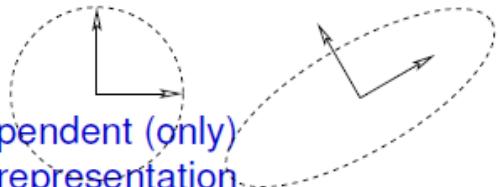
$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mu_w \underbrace{\mathbf{y}_w \mathbf{y}_w^T}_{\text{rank-one}} \quad \text{where } \mu_w = \frac{1}{\sum_{i=1}^{\mu} \mathbf{w}_i^2} \geq 1$$

from [Auger, p. 42]

Rank-One Update: Summary

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

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 \mathbf{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 representation**
components are independent (only) in the new representation
- learns a new (Mahalanobis) metric
variable metric method
- approximates the **inverse Hessian** on quadratic functions
transformation into the sphere function
- for $\mu = 1$: conducts a **natural gradient ascent** on the distribution \mathcal{N}
entirely independent of the given coordinate system

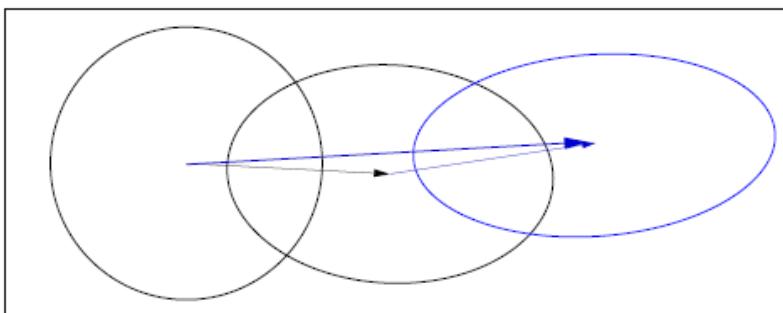
Evolution Path

Cumulation

The Evolution Path

Evolution Path

Conceptually, the evolution path is the **search path** the strategy takes over a **number of generation steps**. It can be expressed as a sum of consecutive *steps* of the mean $\textcolor{green}{m}$.



An exponentially weighted sum of steps y_w is used

$$\textcolor{green}{p}_c \propto \sum_{i=0}^g \underbrace{(1 - \textcolor{blue}{c}_c)^{g-i}}_{\text{exponentially fading weights}} \ y_w^{(i)}$$

The recursive construction of the evolution path (cumulation):

$$\textcolor{green}{p}_c \leftarrow \underbrace{(1 - \textcolor{blue}{c}_c)}_{\text{decay factor}} \textcolor{green}{p}_c + \underbrace{\sqrt{1 - (1 - \textcolor{blue}{c}_c)^2} \sqrt{\mu_w}}_{\text{normalization factor}} \underbrace{y_w}_{\text{input} = \frac{\textcolor{green}{m} - \textcolor{green}{m}_{\text{old}}}{\sigma}}$$

where $\mu_w = \frac{1}{\sum \textcolor{blue}{w}_i^2}$, $c_c \ll 1$. History information is accumulated in the evolution path.

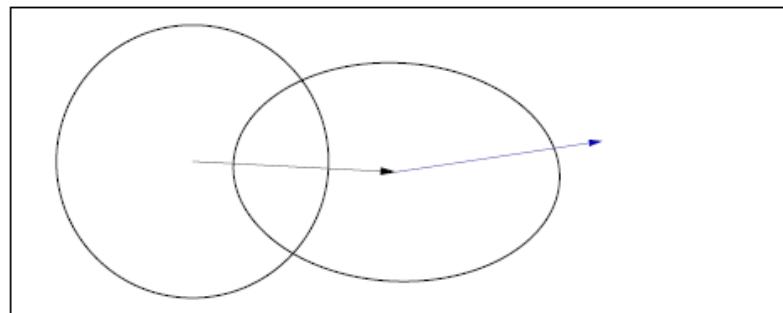
from [Auger, p. 44]

Utilizing the Evolution Path

Cumulation

Utilizing the Evolution Path

We used $y_w y_w^T$ for updating \mathbf{C} . Because $y_w y_w^T = -y_w (-y_w)^T$ the sign of y_w is lost.



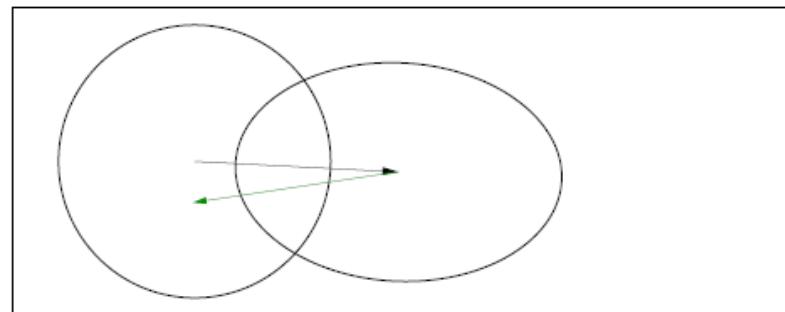
from [Auger, p. 45]

Utilizing the Evolution Path

Cumulation

Utilizing the Evolution Path

We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



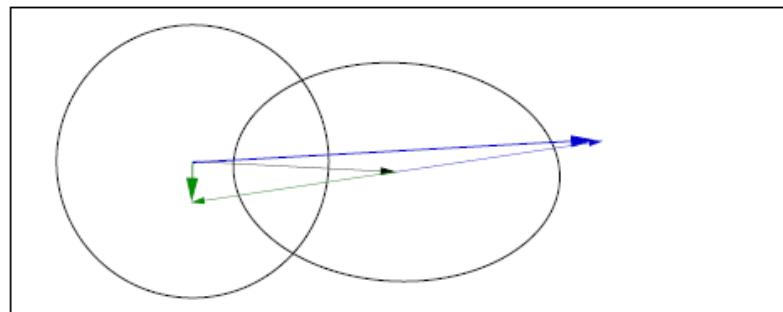
from [Auger, p. 45]

Utilizing the Evolution Path

Cumulation

Utilizing the Evolution Path

We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



The sign information is (re-)introduced by using the *evolution path*.

$$\begin{aligned} \mathbf{p}_c &\leftarrow \underbrace{(1 - c_c)}_{\text{decay factor}} \mathbf{p}_c + \underbrace{\sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w}}_{\text{normalization factor}} \mathbf{y}_w \\ \mathbf{C} &\leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \underbrace{\mathbf{p}_c \mathbf{p}_c^T}_{\text{rank-one}} \end{aligned}$$

where $\mu_w = \sum \mathbf{w}_i^2$, $c_c \ll 1$.

from [Auger, p. 45]

Rank- μ Update

Rank- μ Update

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

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

from [Auger, p. 47]

Rank- μ Update

Rank- μ Update

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

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

The matrix

$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

computes a weighted mean of the outer products of the best μ steps and has rank $\min(\mu, n)$ with probability one.

from [Auger, p. 47]

Rank- μ Update

Rank- μ Update

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

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

The matrix

$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

computes a weighted mean of the outer products of the best μ steps and has rank $\min(\mu, n)$ with probability one.

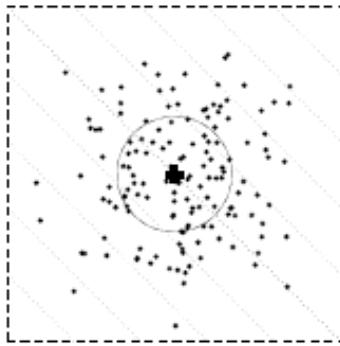
The rank- μ update then reads

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_\mu$$

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

from [Auger, p. 47]

Illustration of Rank- μ Update

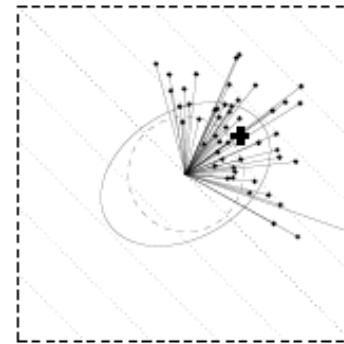
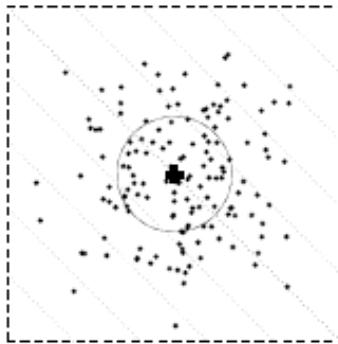


$$x_i = \mathbf{m} + \sigma y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$$

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

from [Auger, p. 48]

Illustration of Rank- μ Update



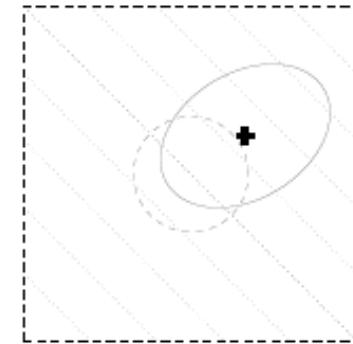
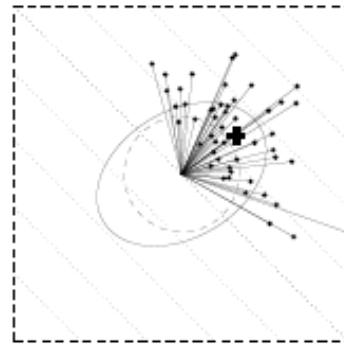
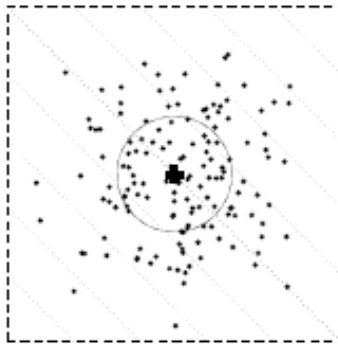
$$x_i = \textcolor{blue}{m} + \textcolor{brown}{\sigma} y_i, \quad y_i \sim \mathcal{N}(0, \textcolor{blue}{C}) \quad C_\mu = \frac{1}{\mu} \sum_{i:\lambda} y_{i:\lambda} y_{i:\lambda}^T \\ \textcolor{blue}{C} \leftarrow (1 - 1) \times \textcolor{blue}{C} + 1 \times C_\mu$$

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

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

from [Auger, p. 48]

Illustration of Rank- μ Update



$$x_i = \mathbf{m} + \sigma y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \mathbf{C}_\mu = \frac{1}{\mu} \sum_{i:\lambda} y_{i:\lambda} y_{i:\lambda}^T \quad \mathbf{C} \leftarrow (1 - 1) \times \mathbf{C} + 1 \times \mathbf{C}_\mu \quad \mathbf{m}_{\text{new}} \leftarrow \mathbf{m} + \frac{1}{\mu} \sum_{i:\lambda}$$

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

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

new distribution

from [Auger, p. 48]

Rank- μ Update: Summary

The rank- μ update

- increases the possible learning rate for large populations
 "large" when $\lambda \geq 3n + 10$
- is the primary mechanism whenever a large population size is used
- can be easily combined with rank-one update

CMA-ES in a Nutshell

The CMA-ES

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

Initialize: $\mathbf{C} = \mathbf{I}$, and $\mathbf{p}_c = \mathbf{0}$,

Set: $c_c \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n$, $c_\mu \approx 1/n$, $w_{i=1\dots\lambda}$ such that $\mu_w = \frac{\sum_{i=1}^\lambda w_i}{n}$,

Promised:

Understand the main principles
of this state-of-the-art algorithm.

While not terminate

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

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

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

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

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

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

Not covered on this slide: termination, restarts, useful output, boundaries and encoding

CMA-ES in a Nutshell

The CMA-ES

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

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

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

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

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

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

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

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

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

Not covered on this slide: termination, restarts, useful output, boundaries and encoding

Strategy Internal Parameters

- related to selection and recombination
 - ▶ λ , offspring number, new solutions sampled, population size
 - ▶ μ , parent number, solutions involved in updates of m , C , and σ
 - ▶ $w_{i=1,\dots,\mu}$, recombination weights
- related to C -update
 - ▶ c_c , decay rate for the evolution path
 - ▶ c_1 , learning rate for rank-one update of C
 - ▶ c_μ , learning rate for rank- μ update of C
- related to σ -update
 - ▶ c_σ , decay rate of the evolution path
 - ▶ d_σ , damping for σ -change

Parameters were identified in carefully chosen experimental set ups. **Parameters do not in the first place depend on the objective function** and are not meant to be in the users choice.

Only(?) the population size λ (and the initial σ) might be reasonably varied in a wide range,
depending on the objective function

Useful: restarts with increasing population size (IPOP)

Experimental Considerations

Experimentum Crucis with CMA-ES

CMA-ES Summary

The Experimentum Crucis

Experimentum Crucis (0)

What did we want to achieve?

- reduce any convex-quadratic function

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

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

to the sphere model

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

without use of derivatives

- lines of equal density align with lines of equal fitness

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

in a stochastic sense

◀ □ ▶ ⏪ ⏩ ⏴ from [Hansen, p. 91]

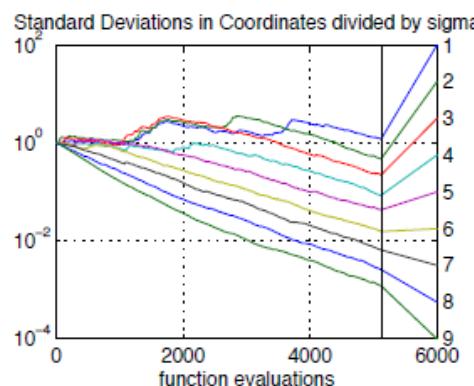
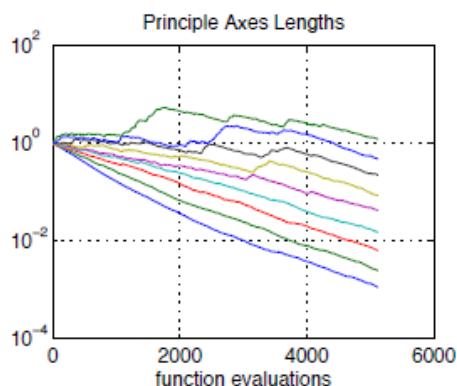
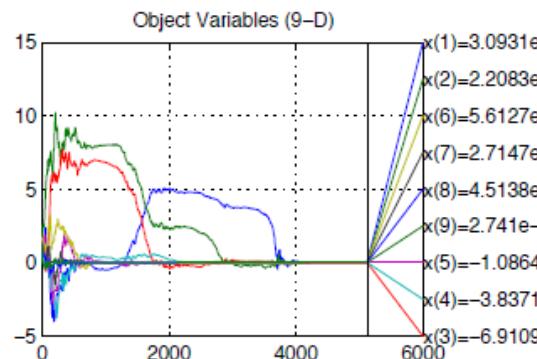
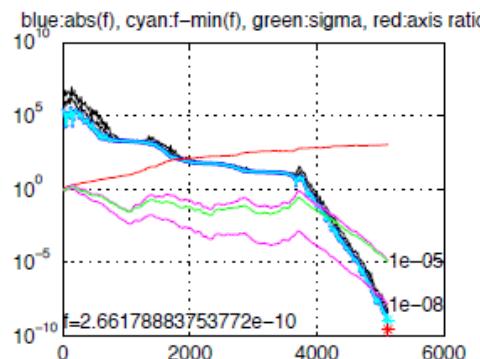
Experimentum Crucis with CMA-ES

CMA-ES Summary

The Experimentum Crucis

Experimentum Crucis (1)

f convex quadratic, separable



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

from [Hansen, p. 92]

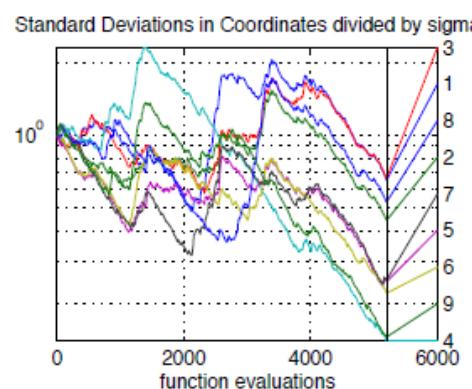
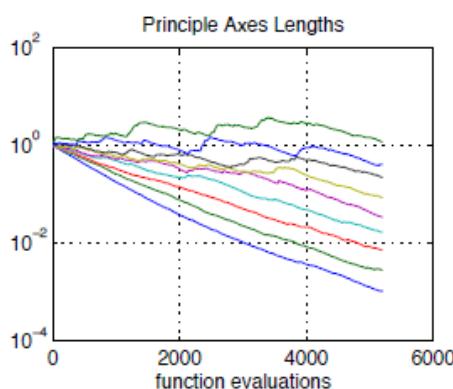
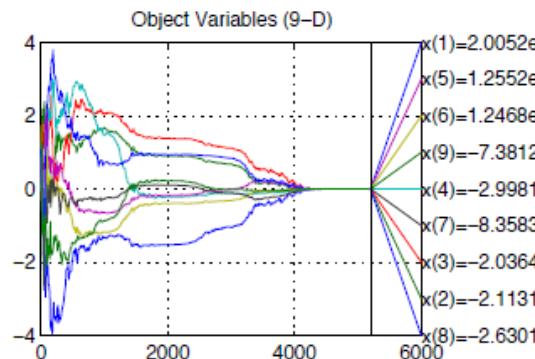
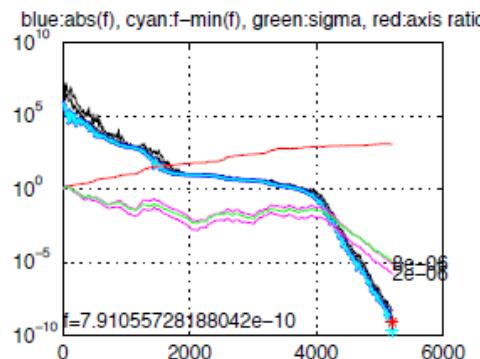
Experimentum Crucis with CMA-ES

CMA-ES Summary

The Experimentum Crucis

Experimentum Crucis (2)

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



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

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

from [Hansen, p. 93]

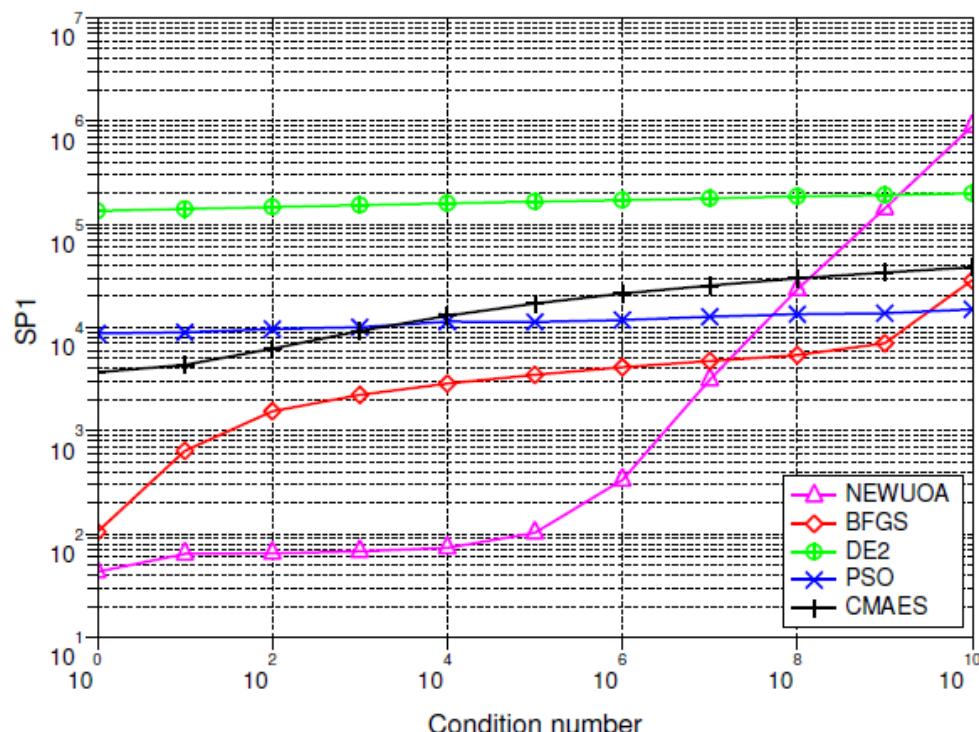
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f convex quadratic, separable with varying condition number α

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



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

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

\mathbf{H} diagonal

g identity (for BFGS and NEWUOA)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁴ to reach the target function value of $g^{-1}(10^{-9})$

¹⁴

Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

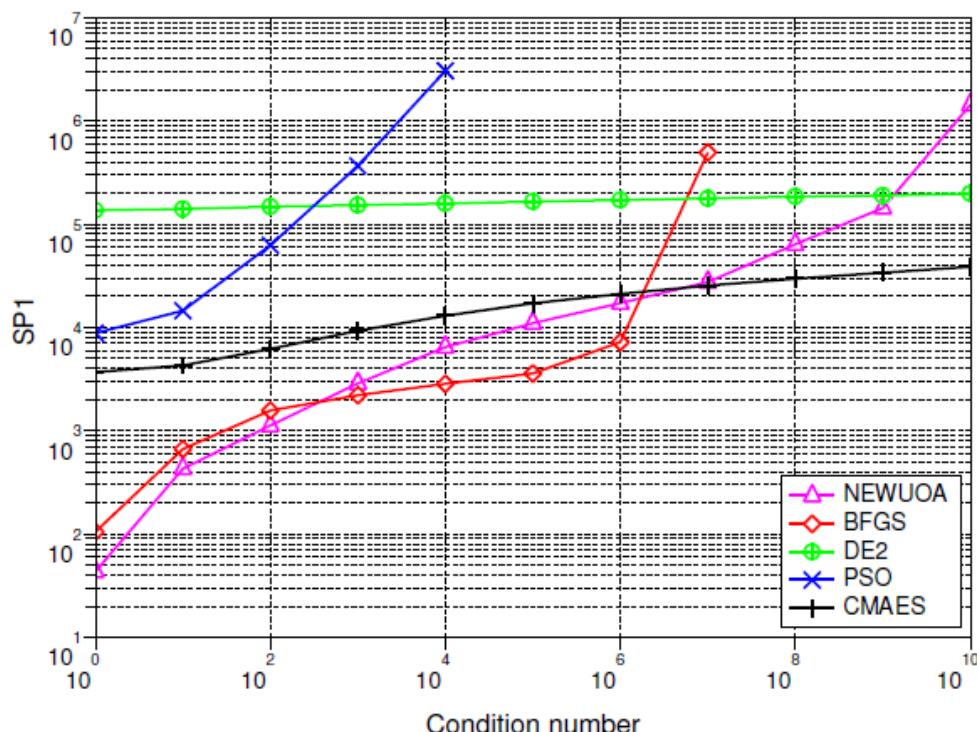
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

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

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



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

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

\mathbf{H} full

g identity (for **BFGS** and **NEWUOA**)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁵ to reach the target function value of $g^{-1}(10^{-9})$

¹⁵

Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

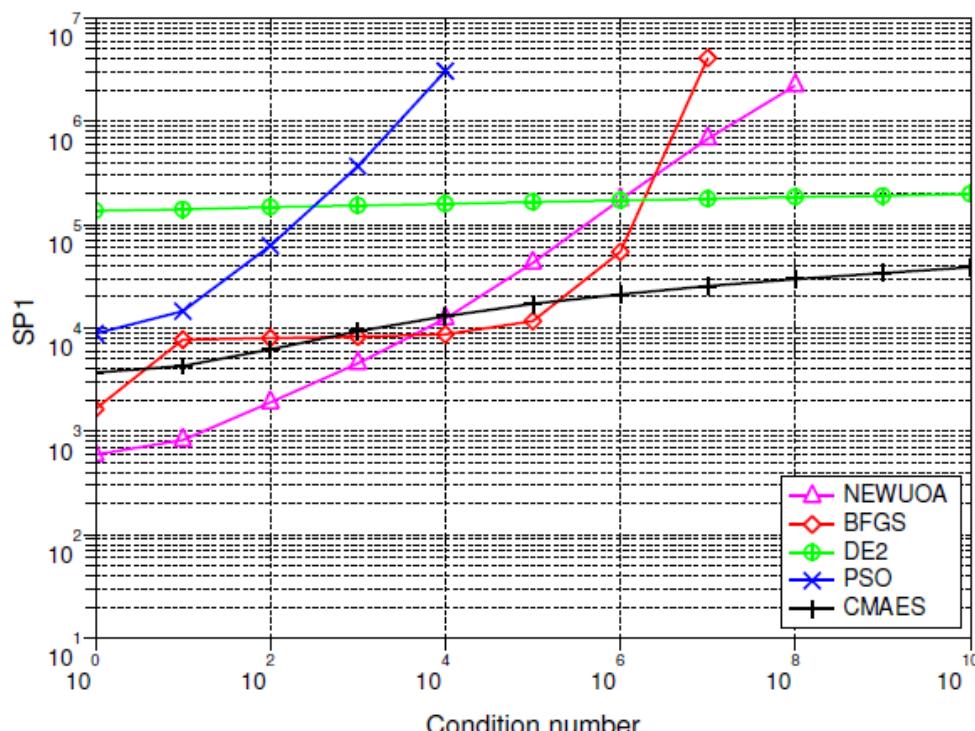
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

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

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



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

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

\mathbf{H} full

$g : x \mapsto x^{1/4}$ (for BFGS and NEWUOA)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁶ to reach the target function value of $g^{-1}(10^{-9})$

¹⁶

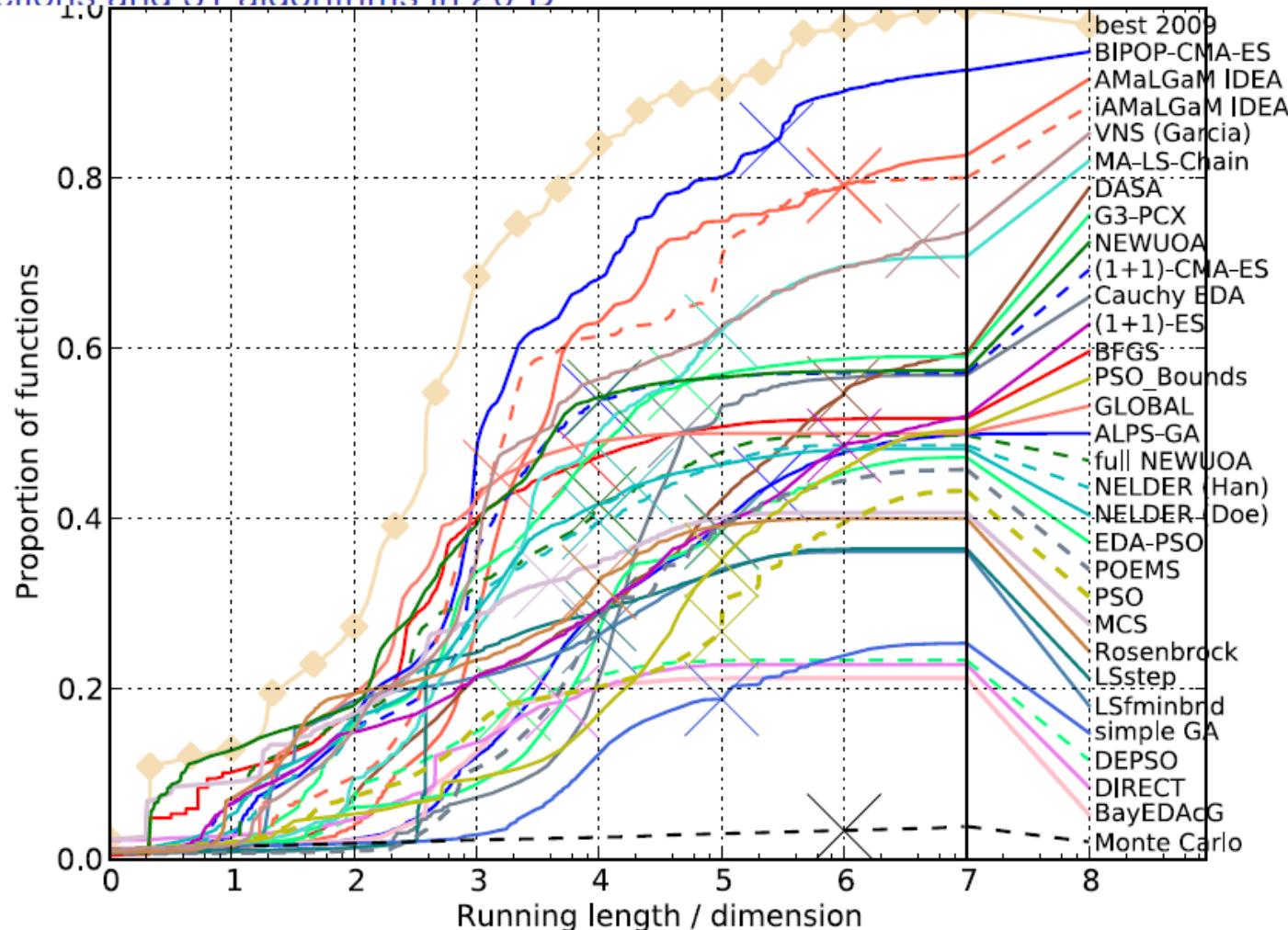
Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

Performance on BBOB Testbed: Data Profile

Comparing Experiments

Comparison during BBOB at GECCO 2009

24 functions and 31 algorithms in 20-D



Main Characteristics of (CMA) Evolution Strategies

- ① Multivariate normal distribution to generate new search points
follows the maximum entropy principle
- ② Rank-based selection
implies invariance, same performance on $g(f(x))$ for any increasing g
more invariance properties are featured
- ③ Step-size control facilitates fast (log-linear) convergence and
possibly linear scaling with the dimension
in CMA-ES based on an **evolution path** (a non-local trajectory)
- ④ Covariance matrix adaptation (CMA) **increases the likelihood of**
previously successful steps and can improve performance by
orders of magnitude
 - the update follows the natural gradient
 - $\mathbf{C} \propto \mathbf{H}^{-1} \iff$ adapts a variable metric
 - \iff new (rotated) problem representation
 - $\implies f : \mathbf{x} \mapsto g(\mathbf{x}^T \mathbf{H} \mathbf{x})$ reduces to $\mathbf{x} \mapsto \mathbf{x}^T \mathbf{x}$

Limitations

of CMA Evolution Strategies

- internal CPU-time: $10^{-8}n^2$ seconds per function evaluation on a 2GHz PC, tweaks are available
1 000 000 f -evaluations in 100-D take 100 seconds *internal CPU-time*
- better methods are presumably available in case of
 - ▶ partly separable problems
 - ▶ specific problems, for example with cheap gradients
specific methods
 - ▶ small dimension ($n \ll 10$)
for example Nelder-Mead
 - ▶ small running times (number of f -evaluations $< 100n$)
model-based methods

Conclusions

I hope it became clear...

...that CMA-ES samples according to multivariate normal distributions
...how CMA-ES updates its mean, stepsize, and covariance matrix
...and what are the invariance properties of CMA-ES