

# Introduction to Optimization

## Derivative-Free Optimization I: CMA-ES

December 12, 2016  
École Centrale Paris, Châtenay-Malabry, France



Dimo Brockhoff  
Inria Saclay – Ile-de-France

# Course Overview

Date		Topic
Fri, 7.10.2016		Introduction
Fri, 28.10.2016	D	Introduction to Discrete Optimization + Greedy algorithms I
Fri, 4.11.2016	D	Greedy algorithms II + Branch and bound
Fri, 18.11.2016	D	Dynamic programming
Mon, 21.11.2016 in S103-S105	D	Approximation algorithms <del>and heuristics</del>
Fri, 25.11.2016 in S103-S105	C	Randomized Search Heuristics + Intro. to Continuous Opt. I
Mon, 28.11.2016 in S103-S105	C	Introduction to Continuous Optimization II
Mon, 5.12.2016 in S103-S105	C	Introduction to Continuous Optimization III
Fri, 9.12.2016	C	Constrained Optimization + Descent Methods
Mon, 12.12.2016 in S103-S105	C	Derivative Free Optimization I: CMA-ES
Fri, 16.12.2016	C	Derivative Free Optimization II: Benchmarking Optimizers with the COCO platform
Wed, 4.1.2017		Exam

if not indicated otherwise, classes take place in S115-S117

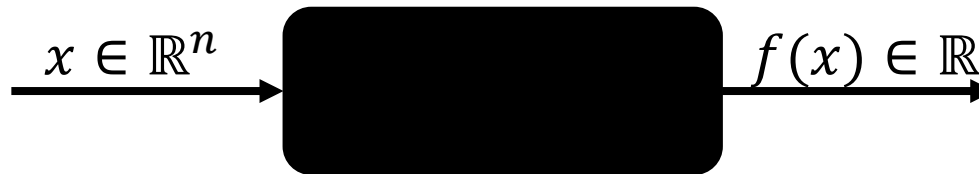
# Solution to the Exercise: Comparing Gradient-Based Algorithms on Convex Quadratic Functions

`http://researchers.lille.inria.fr/  
~brockhof/introoptimization/`

# Derivative-Free Optimization

# Derivative-Free Optimization (DFO)

DFO = blackbox optimization



## Why blackbox scenario?

- gradients are not always available (binary code, no analytical model, ...)
- or not useful (noise, non-smooth, ...)
- problem domain specific knowledge is used only within the black box, e.g. within an appropriate encoding
- some algorithms are furthermore function-value-free, i.e. *invariant* wrt. monotonous transformations of  $f$ .

# Derivative-Free Optimization Algorithms

- (gradient-based algorithms which approximate the gradient by finite differences)
- coordinate descent
- **pattern search** methods, e.g. Nelder-Mead
- surrogate-assisted algorithms, e.g. NEWUOA or other **trust-region methods**
- other **function-value-free algorithms**
  - typically stochastic
  - evolution strategies (ESs) and Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
  - differential evolution
  - particle swarm optimization
  - simulated annealing
  - ...

# Derivative-Free Optimization Algorithms

- (gradient-based algorithms which approximate the gradient by finite differences)
- coordinate descent
- **pattern search** methods, e.g. **Nelder-Mead**
- surrogate-assisted algorithms, e.g. NEWUOA or other **trust-region methods**
- other **function-value-free algorithms**
  - typically stochastic
  - evolution strategies (ESs) and Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
  - differential evolution
  - particle swarm optimization
  - simulated annealing
  - ...

# Downhill Simplex Method by Nelder and Mead

While not happy do:

[assuming minimization of  $f$  and that  $x_1, \dots, x_{n+1} \in \mathbb{R}^n$  form a simplex]

**1) Order** according to the values at the vertices:  $f(x_1) \leq f(x_2) \leq \dots \leq f(x_{n+1})$

**2)** Calculate  $x_o$ , the centroid of all points except  $x_{n+1}$ .

**3) Reflection**

Compute reflected point  $x_r = x_o + \alpha (x_o - x_{n+1})$  ( $\alpha > 0$ )

If  $x_r$  better than second worst, but not better than best:  $x_{n+1} := x_r$ , and go to 1)

**4) Expansion**

If  $x_r$  is the best point so far: compute the expanded point

$$x_e = x_o + \gamma (x_r - x_o) (\gamma > 0)$$

If  $x_e$  better than  $x_r$  then  $x_{n+1} := x_e$  and go to 1)

Else  $x_{n+1} := x_r$  and go to 1)

Else (i.e. reflected point is not better than second worst) continue with 5)

**5) Contraction** (here:  $f(x_r) \geq f(x_n)$ )

Compute contracted point  $x_c = x_o + \rho (x_{n+1} - x_o)$  ( $0 < \rho \leq 0.5$ )

If  $f(x_c) < f(x_{n+1})$ :  $x_{n+1} := x_c$  and go to 1)

Else go to 6)

**6) Shrink**

$x_i = x_1 + \sigma (x_i - x_1)$  for all  $i \in \{2, \dots, n+1\}$  and go to 1)

*J. A Nelder and R. Mead (1965). "A simplex method for function minimization".  
Computer Journal. 7: 308–313. doi:10.1093/comjnl/7.4.308*



# Downhill Simplex Method by Nelder and Mead

- initial simplex is important: hence restarts necessary to have good performance
- illustration of working principles at [https://en.wikipedia.org/wiki/Nelder%E2%80%93Mead\\_method](https://en.wikipedia.org/wiki/Nelder%E2%80%93Mead_method)
- nice-to-read paper about the (historical) background of the method: [http://www.math.uiuc.edu/documenta/vol-ismp/42\\_wright-margaret.pdf](http://www.math.uiuc.edu/documenta/vol-ismp/42_wright-margaret.pdf)
- turns out to be quite good in low-dimensional problems (with 2 or 3 variables), but not in high dimension (see also this Friday's exercise)

# Derivative-Free Optimization Algorithms

- (gradient-based algorithms which approximate the gradient by finite differences)
- coordinate descent
- **pattern search** methods, e.g. Nelder-Mead
- surrogate-assisted algorithms, e.g. NEWUOA or other **trust-region methods**
- other **function-value-free algorithms**
  - typically stochastic
  - evolution strategies (ESs) and **Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
  - differential evolution
  - particle swarm optimization
  - simulated annealing
  - ...

# Stochastic Search Template

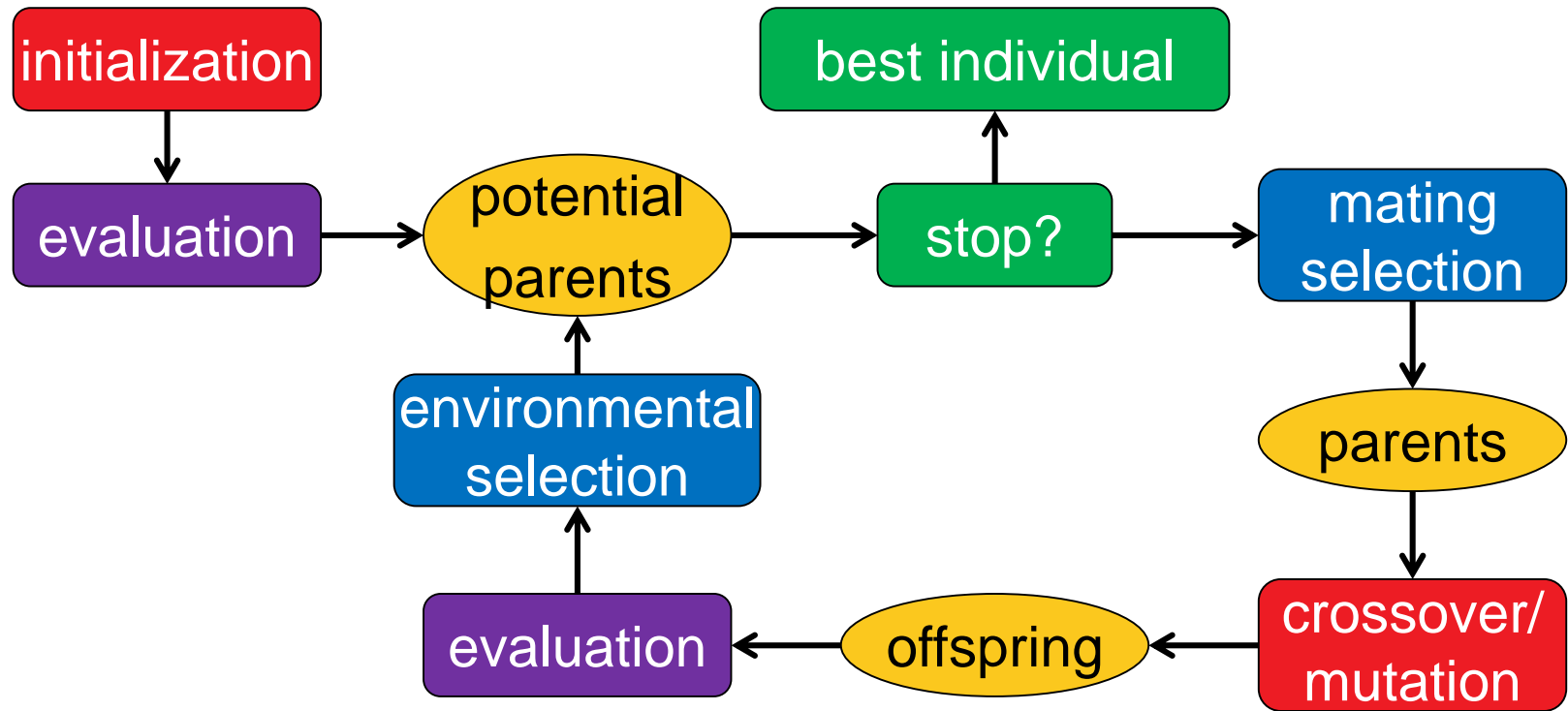
**A stochastic blackbox search template to minimize  $f: \mathbb{R}^n \rightarrow \mathbb{R}$**

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

While happy do:

- Sample distribution  $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$
  - Evaluate  $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$  on  $f$
  - Update parameters  $\theta \leftarrow F_\theta(\theta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$
- 
- All depends on the choice of  $P$  and  $F_\theta$ 
    - deterministic algorithms are covered as well*
  - In Evolutionary Algorithms,  $P$  and  $F_\theta$  are often defined implicitly via their operators.

# Generic Framework of an EA



stochastic operators

“Darwinism”

stopping criteria

Nothing else: just interpretation change

## The CMA-ES

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

**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**

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

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

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

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

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

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

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

## The CMA-ES

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

**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**

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

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

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

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^\top + c_\mu \sum_{i=1}^\mu \mathbf{y}_i \mathbf{y}_i^\top$  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:**  $\mathbf{m} \in \mathbb{R}^n, \sigma \in \mathbb{R}_+, \lambda$

**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**

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

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

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

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^\mu \mathbf{y}_i \mathbf{y}_i^T$  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  $\theta$ , set population size  $\lambda \in \mathbb{N}$

While happy do:

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

For CMA-ES and evolution strategies in general:

sample distributions = multivariate Gaussian distributions

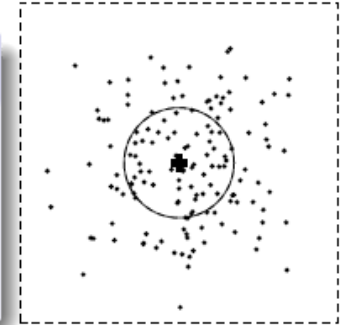
# Sampling New Candidate Solutions (Offspring)

## Evolution Strategies

New search points are sampled normally distributed

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

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
- 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

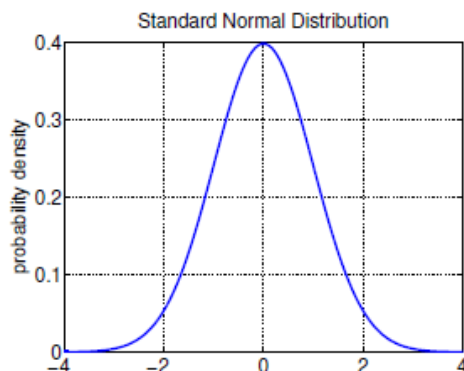
it remains to show how to adapt the parameters, but for now: normal distributions

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  $\mathcal{N}(m, \sigma^2)$

(expected value, variance) =  $(m, \sigma^2)$

density:  $p_{m,\sigma}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right)$

- ▶ 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  $m + \sigma\mathcal{N}(0, 1) = \mathcal{N}(m, \sigma^2)$

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**  $\sigma^2$ .

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

### Covariance Matrix

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

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

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

$$\Sigma = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$$

$\Sigma$  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  $\mathbf{C}$ .

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

from [Auger, p. 13]

# Excursion: Normal Distributions

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

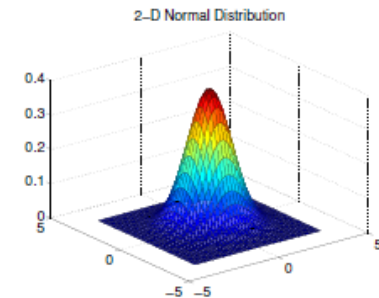
Any multi-variate normal distribution  $\mathcal{N}(\mathbf{m}, \mathbf{C})$  is uniquely determined by its mean value  $\mathbf{m} \in \mathbb{R}^n$  and its symmetric positive definite  $n \times n$  covariance matrix  $\mathbf{C}$ .

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

The mean value  $\mathbf{m}$

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

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



from [Auger, p. 13]

# Excursion: Normal Distributions

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

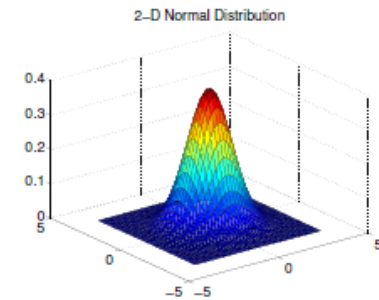
Any multi-variate normal distribution  $\mathcal{N}(\mathbf{m}, \mathbf{C})$  is uniquely determined by its mean value  $\mathbf{m} \in \mathbb{R}^n$  and its symmetric positive definite  $n \times n$  covariance matrix  $\mathbf{C}$ .

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

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

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

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



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

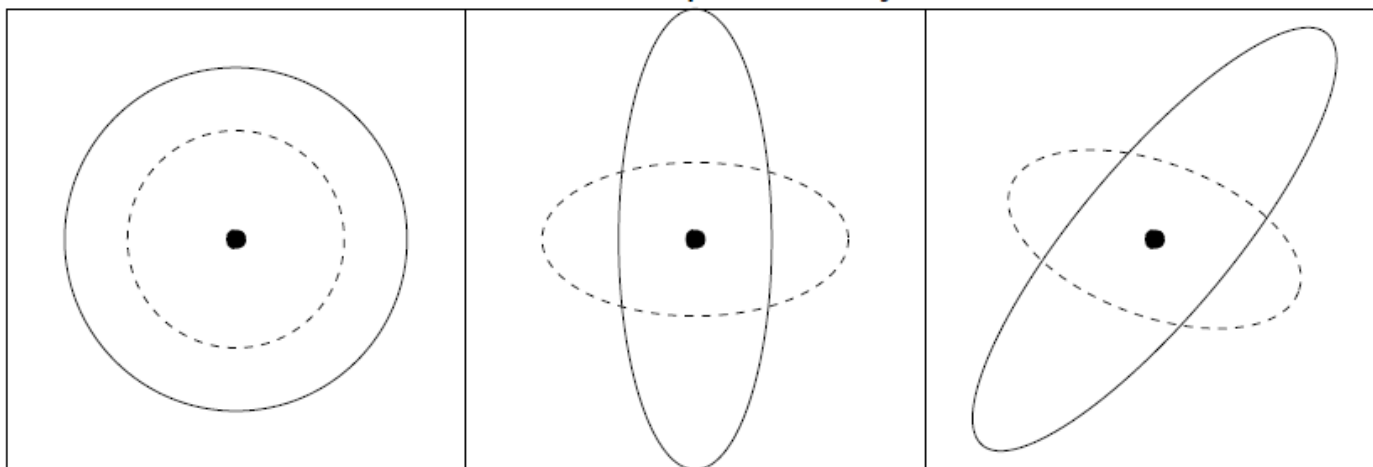
- ▶ determines the shape
- ▶ **geometrical interpretation**: any covariance matrix can be uniquely identified with the iso-density ellipsoid  $\{\mathbf{x} \in \mathbb{R}^n \mid (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}) = 1\}$

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  $\sigma$

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  $\mathbf{A}$ .

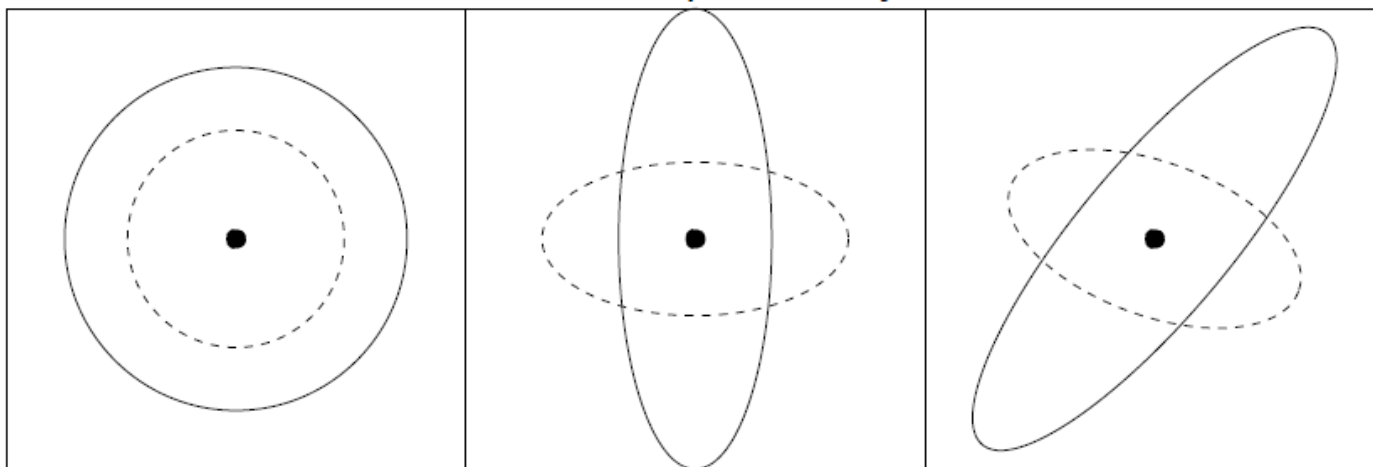
from [Auger, p. 14]



# 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  $\sigma$

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

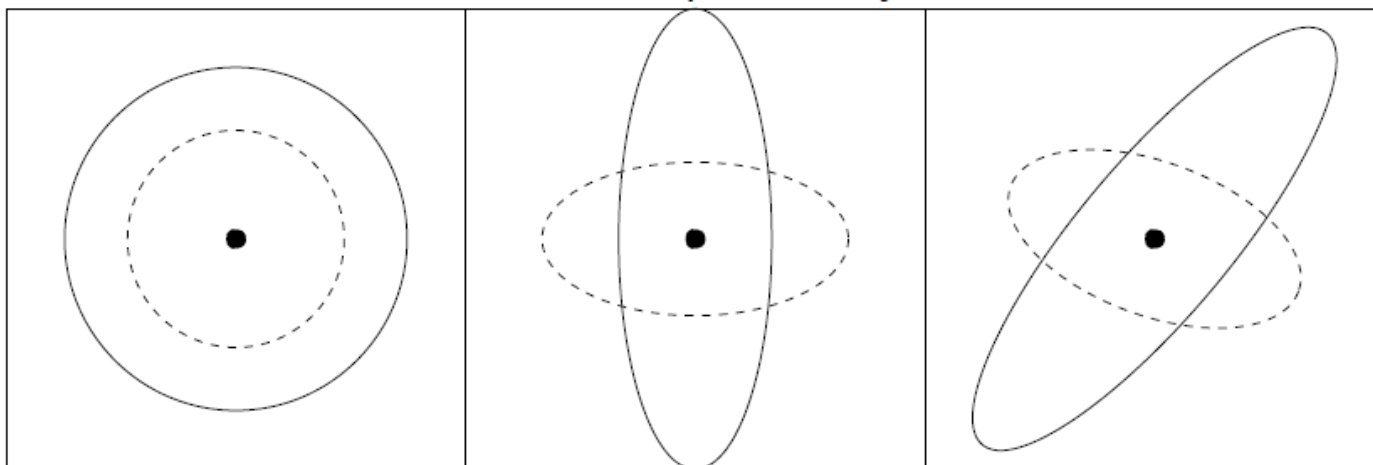
where  $\mathbf{I}$  is the identity matrix (isotropic case) and  $\mathbf{D}$  is a diagonal matrix (reasonable for separable problems) and  $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T)$  holds for all  $\mathbf{A}$ .

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  $\sigma$   
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  
particularily important for **ill-conditioned** problems

$\mathbf{C} \propto \mathbf{H}^{-1}$  on convex quadratic functions

from [Auger, p. 16]

# Adaptation of the Mean

## Evolution Strategies

### Terminology

$\mu$ : # of parents,  $\lambda$ : # of offspring

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

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

$(\mu, \lambda)$ -ES: selection in  $\{\text{offspring}\}$

### $(1 + 1)$ -ES

Sample one offspring from parent  $m$

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

If  $x$  better than  $m$  select

$$m \leftarrow x$$

## 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

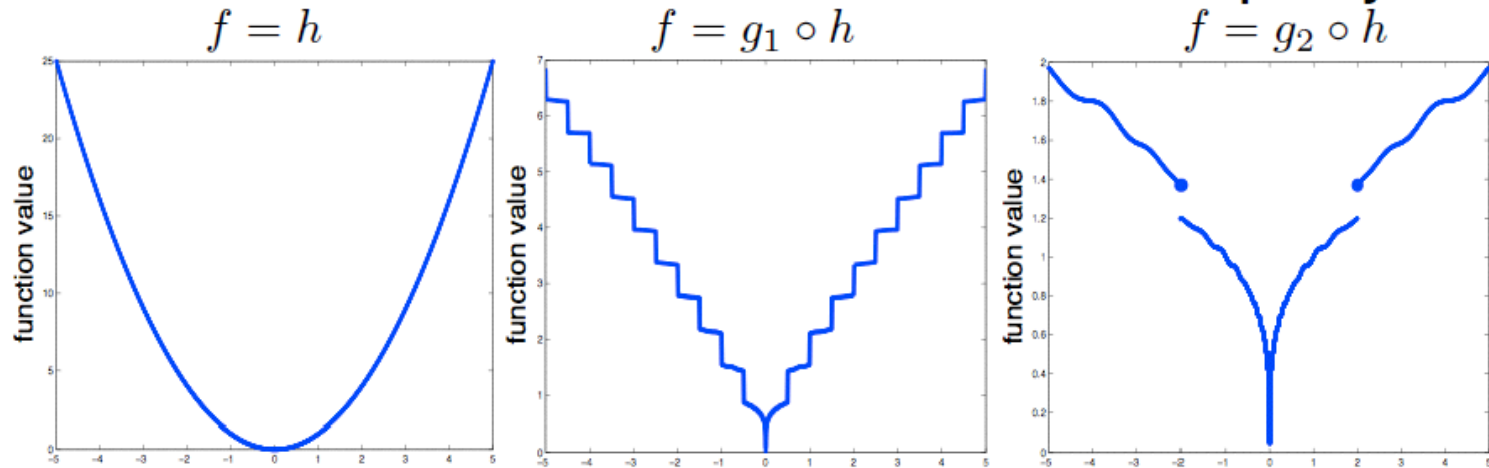
$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$

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

from [Hansen, p. 34]

# 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

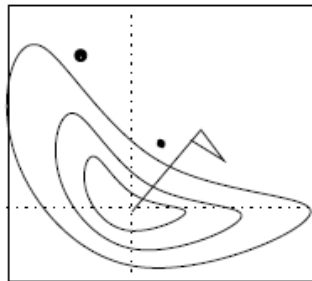
- observations meaningful as a rigorous notion of generalization
- algorithms predictable and/or "robust"

from [Hansen, p. 37]

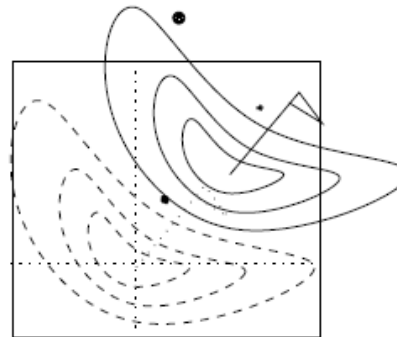
## Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



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



Identical behavior on  $f$  and  $f_a$

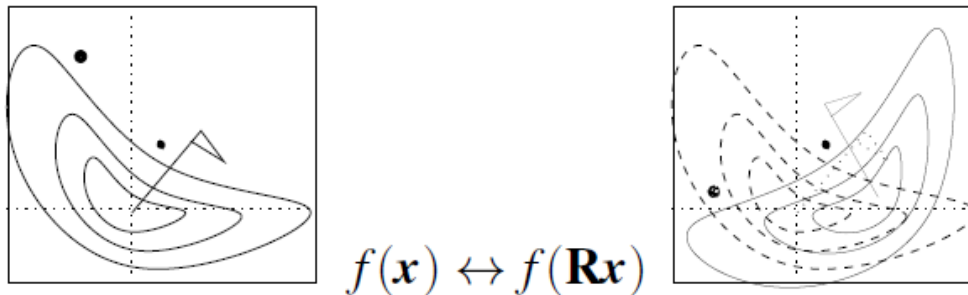
$$\begin{aligned} f &: \mathbf{x} \mapsto f(\mathbf{x}), & \mathbf{x}^{(t=0)} &= \mathbf{x}_0 \\ f_a &: \mathbf{x} \mapsto f(\mathbf{x} - \mathbf{a}), & \mathbf{x}^{(t=0)} &= \mathbf{x}_0 + \mathbf{a} \end{aligned}$$

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



## Rotational Invariance in Search Space

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



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

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

45

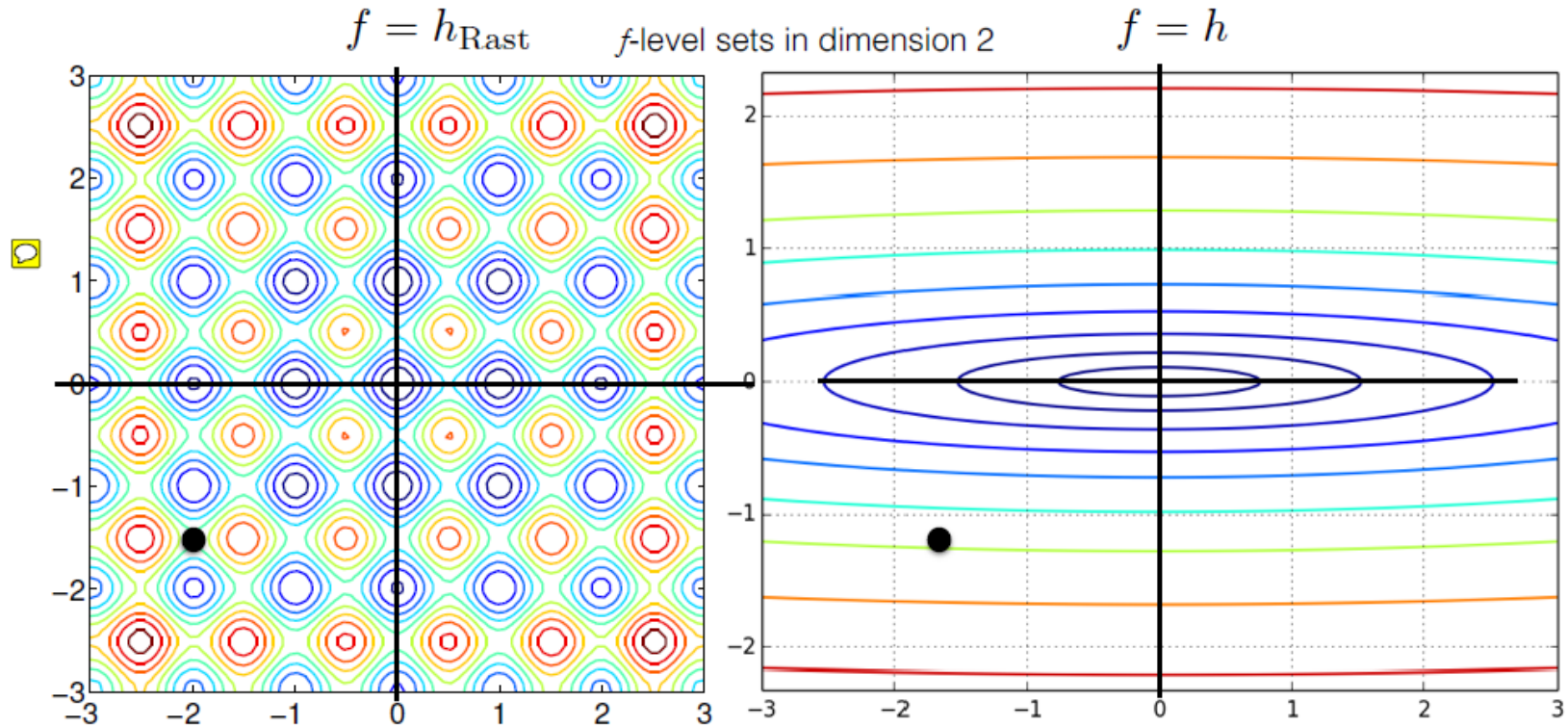
No difference can be observed w.r.t. the argument of  $f$

<sup>4</sup>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

<sup>5</sup>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)

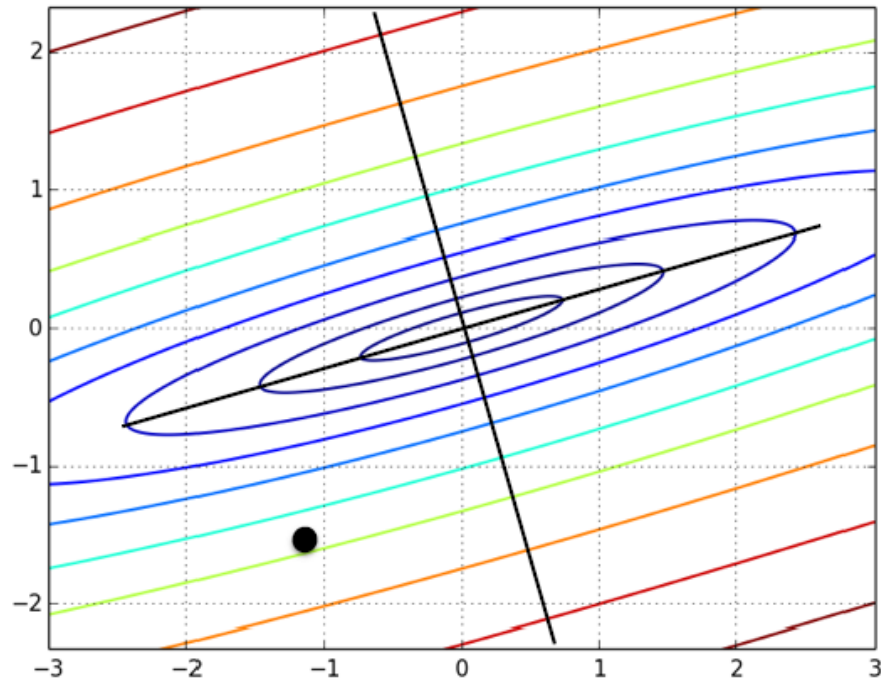
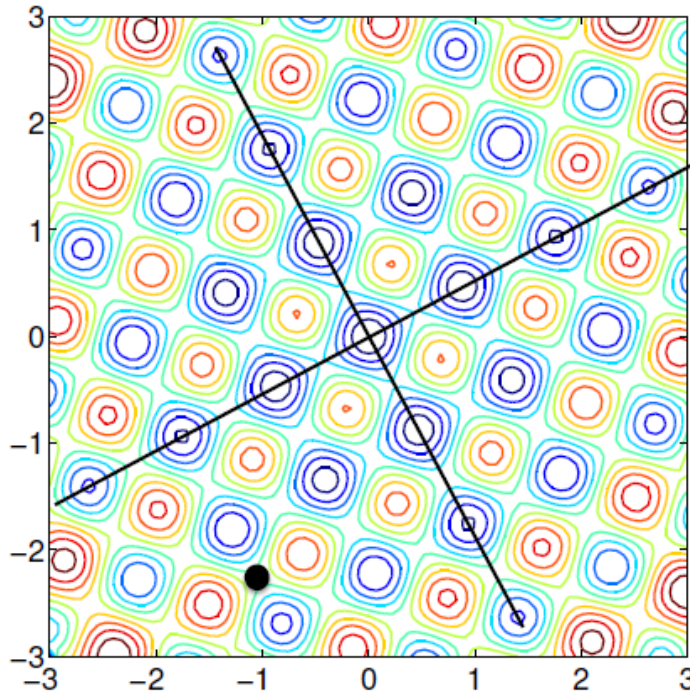
from [Hansen, p. 40

## Invariance Under Rigid Search Space Transformations

$$f = h_{\text{Rast}} \circ R$$

$f$ -level sets in dimension 2

$$f = h \circ R$$



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

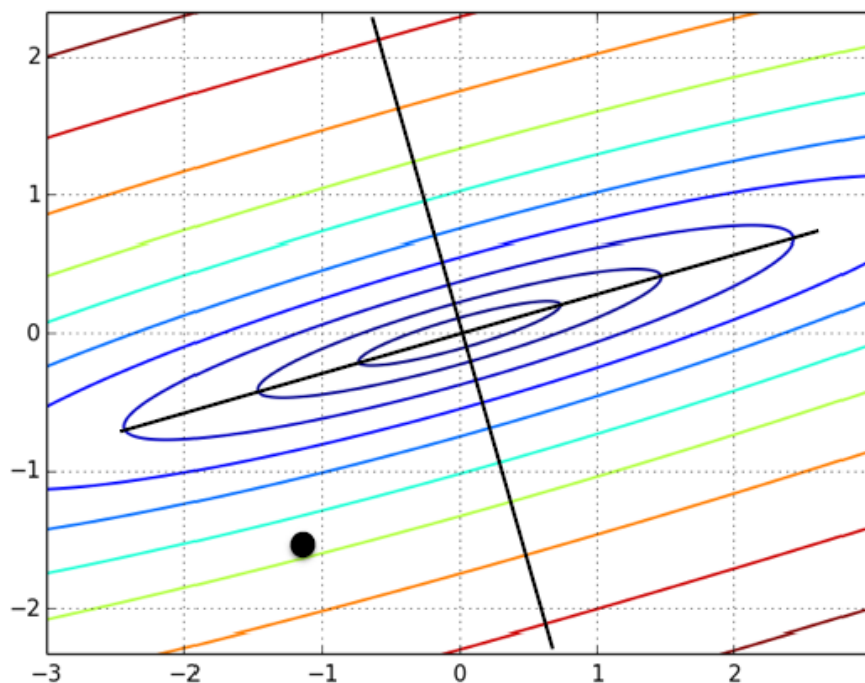
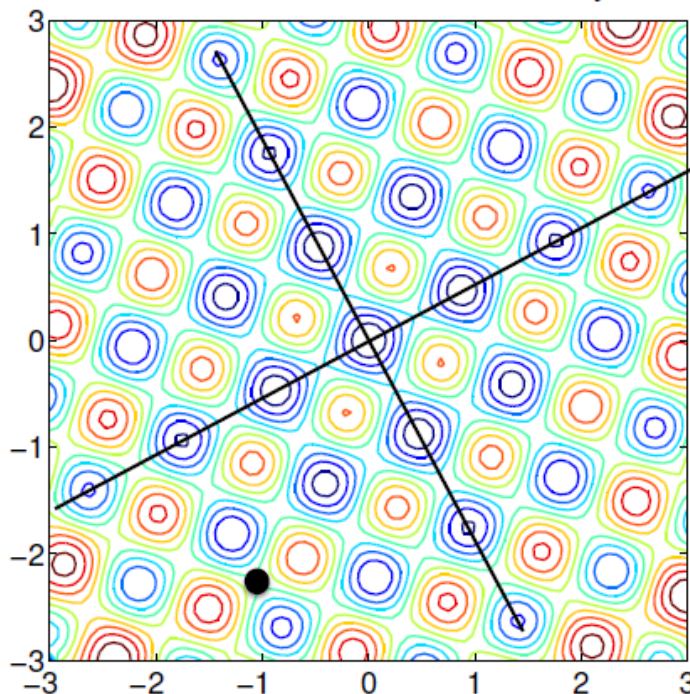
from [Hansen, p. 41]

## Invariance Under Rigid Search Space Transformations

$$f = h_{\text{Rast}} \circ R$$

$f$ -level sets in dimension 2

$$f = h \circ R$$



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

mainly Nelder-Mead and CMA-ES  
have this property

## 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

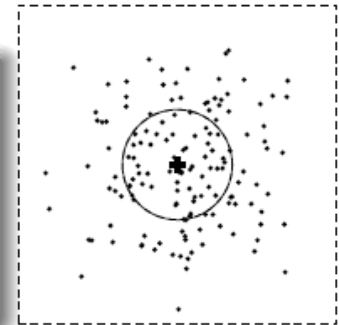
## Evolution Strategies

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



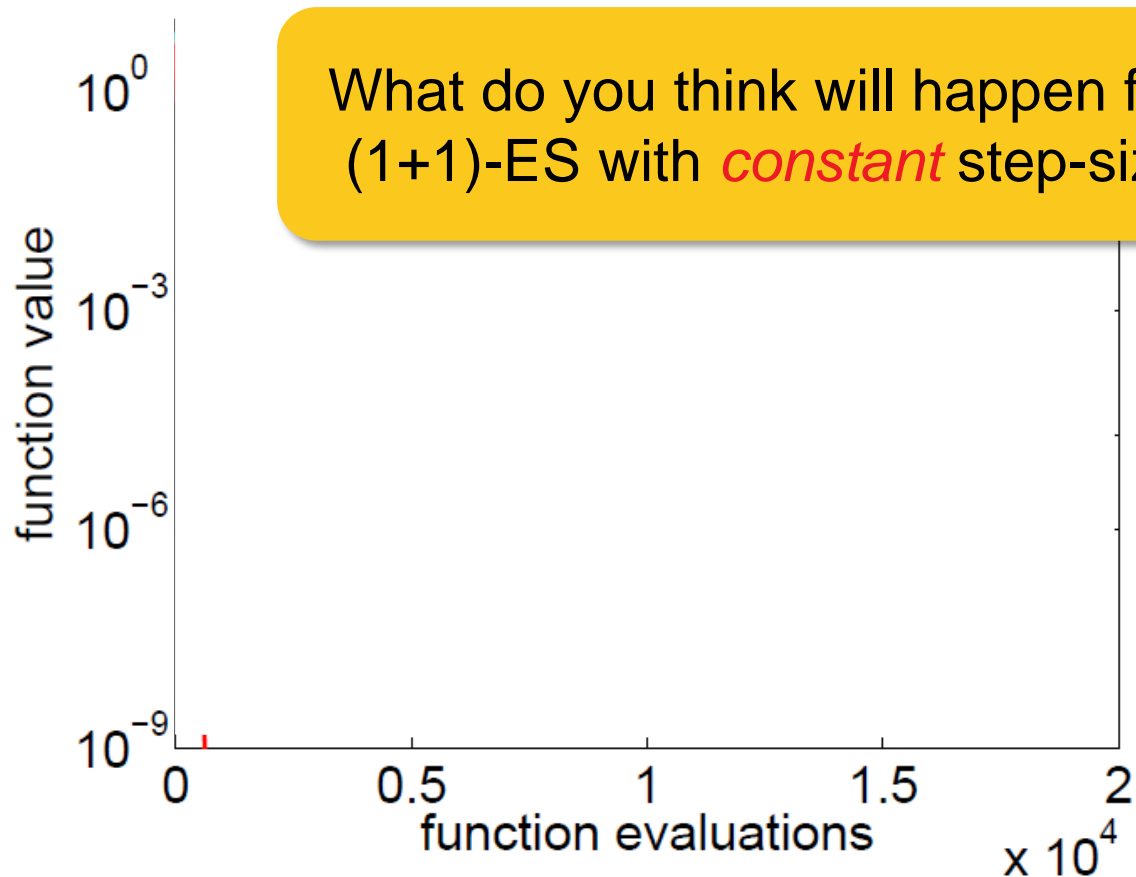
where

- 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  $\sigma$  and  $\mathbf{C}$ .

# Why At All Step-Size Adaptation?

## Why Step-Size Control?



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

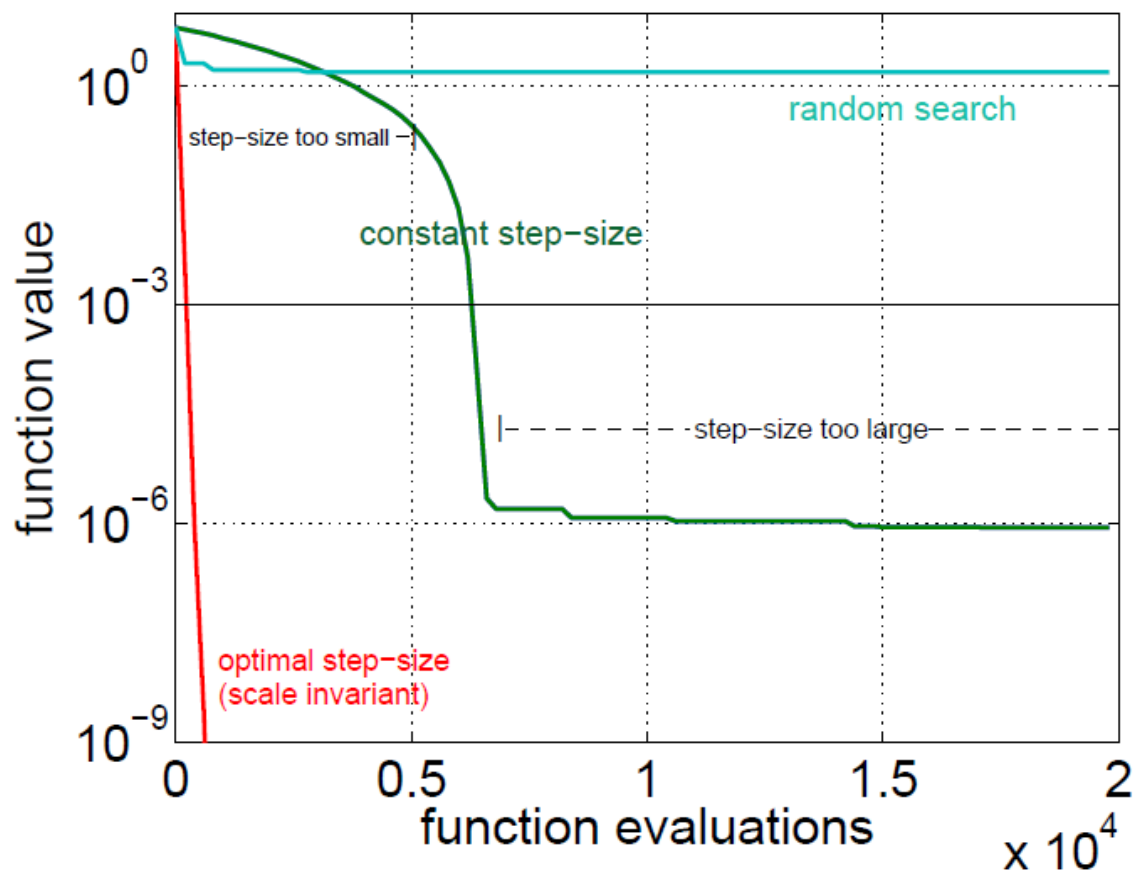
in  $[-0.2, 0.8]^n$   
for  $n = 10$

from [Auger, p. 22]



# Why Step-Size Adaptation?

## Why Step-Size Control?



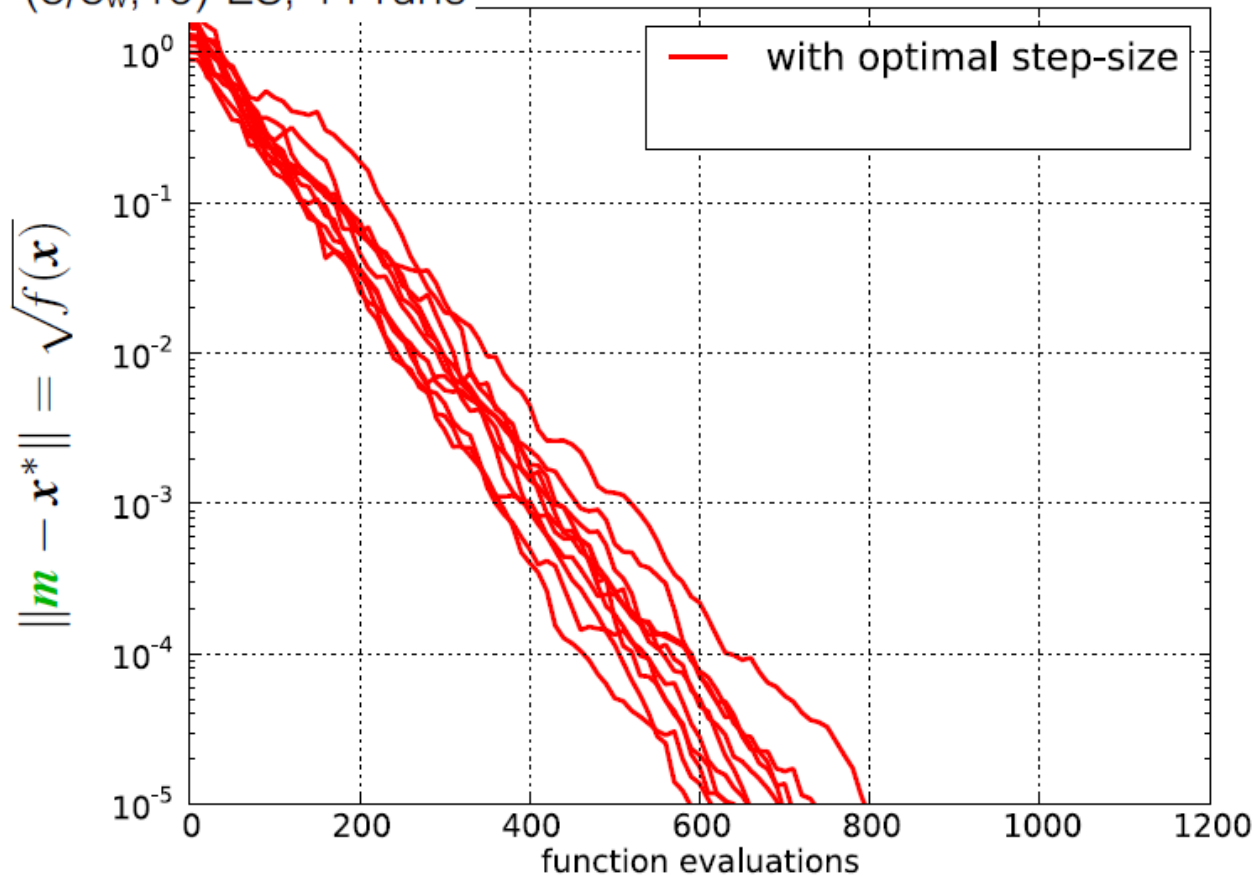
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

from [Auger, p. 22]

## Why Step-Size Control?

(5/5<sub>w</sub>,10)-ES, 11 runs



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

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

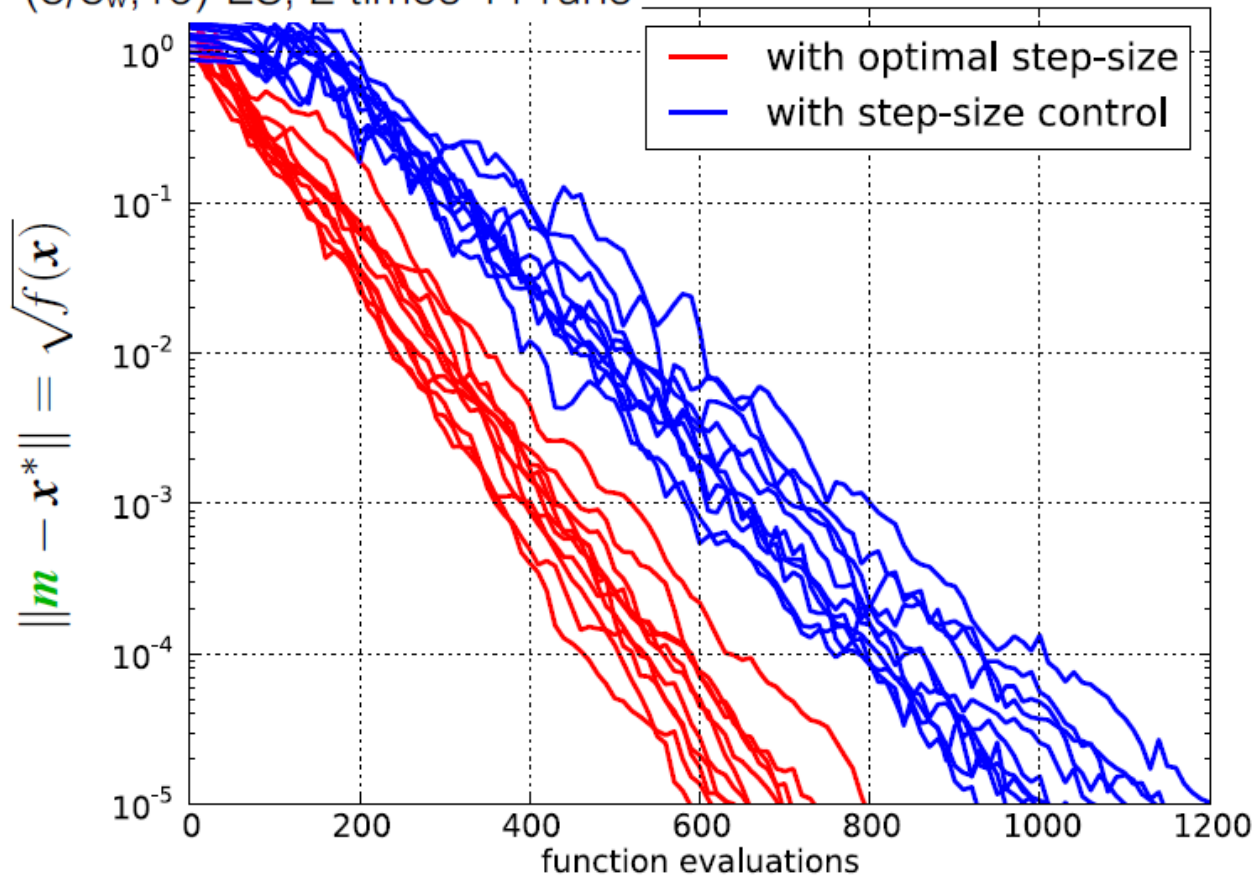
with optimal step-size  $\sigma$

from [Hansen, p. 47]

# Optimal Step-Size vs. Step-Size Control

## Why Step-Size Control?

(5/5<sub>w</sub>,10)-ES, 2 times 11 runs



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

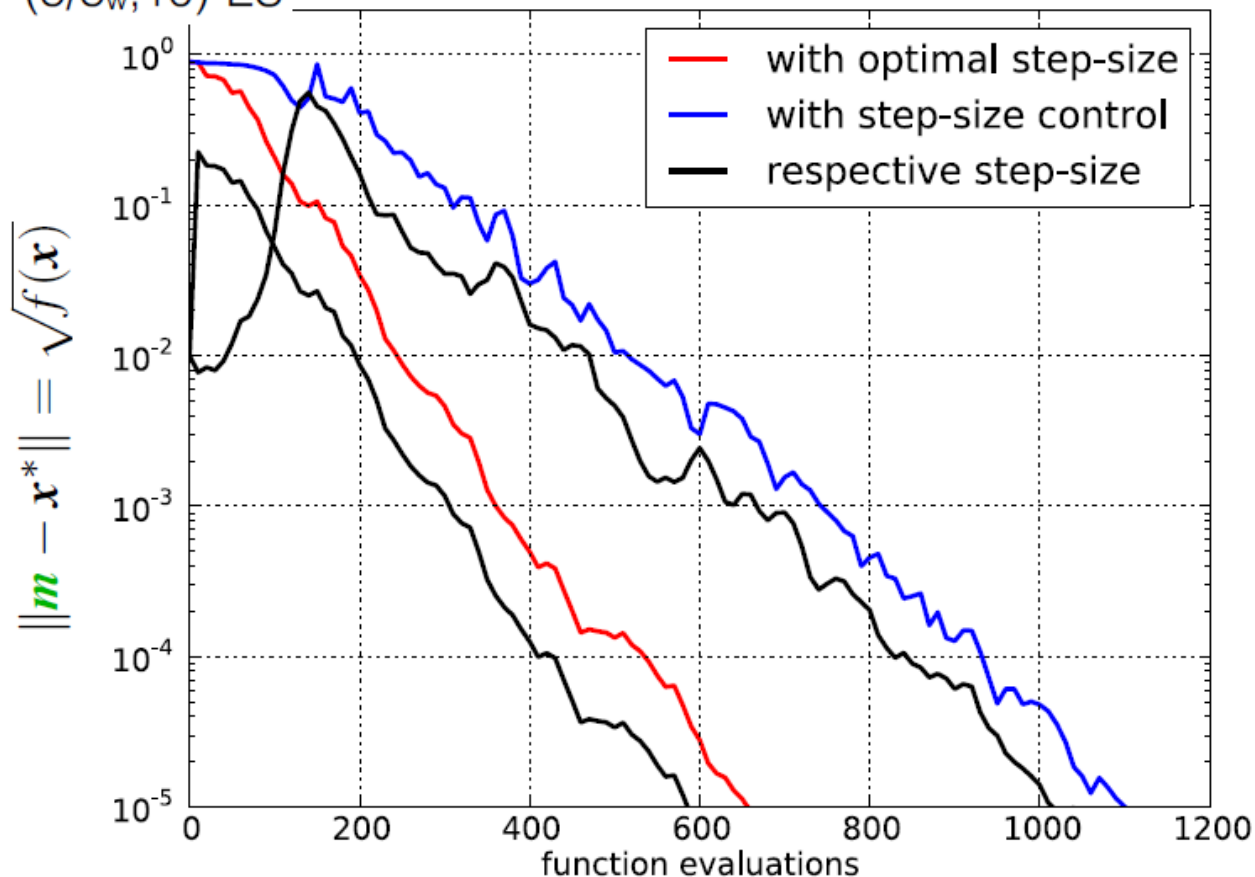
for  $n = 10$  and  
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

with **optimal** versus **adaptive** step-size  $\sigma$  with too small initial  $\sigma$

# Optimal Step-Size vs. Step-Size Control

## Why Step-Size Control?

(5/5<sub>w</sub>, 10)-ES



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

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

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

from [Hansen, p. 49]

# Adapting the Step-Size

## Question:

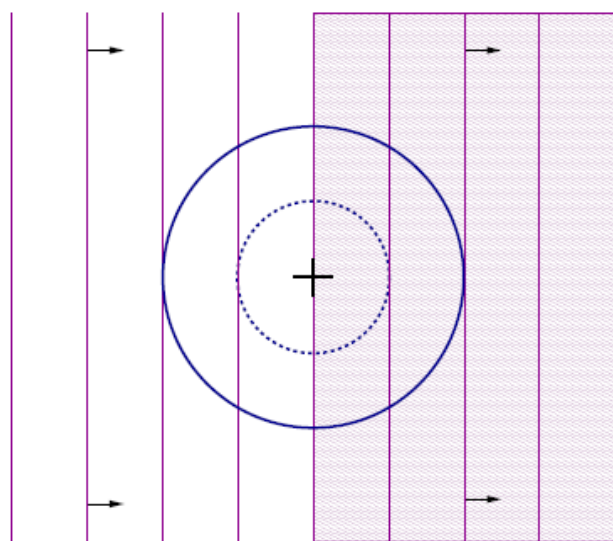
How to actually adapt the step-size during the optimization?

## Most common:

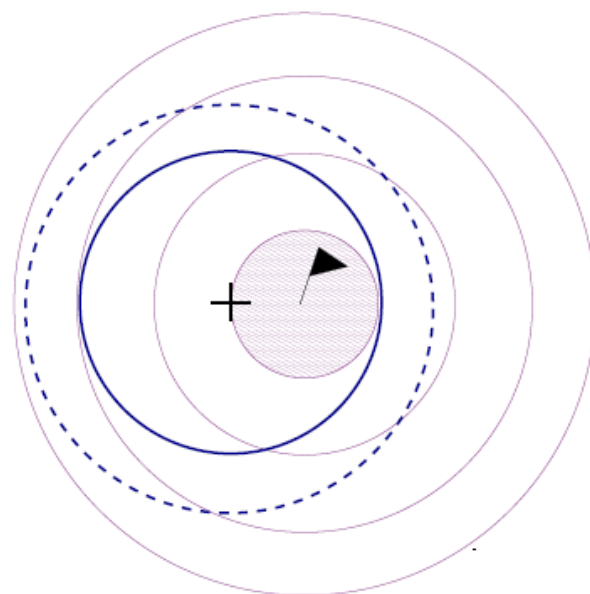
- 1/5 success rule
- Cumulative Step-Size Adaptation (CSA, as in standard CMA-ES)
- others possible (Two-Point Adaptation, self-adaptive step-size, ...)

# One-Fifth Success Rule

## One-fifth success rule



↓  
increase  $\sigma$

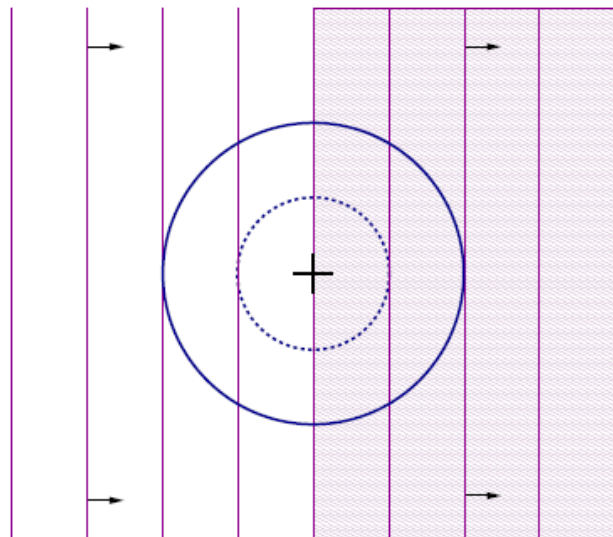


↓  
decrease  $\sigma$

from [Auger, p. 32]

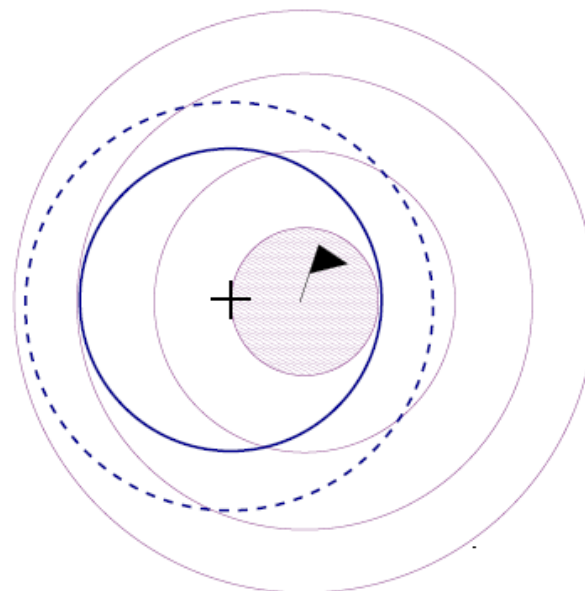
# One-Fifth Success Rule

## 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  $\sigma$  if  $p_s > p_{\text{target}}$   
Decrease  $\sigma$  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

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

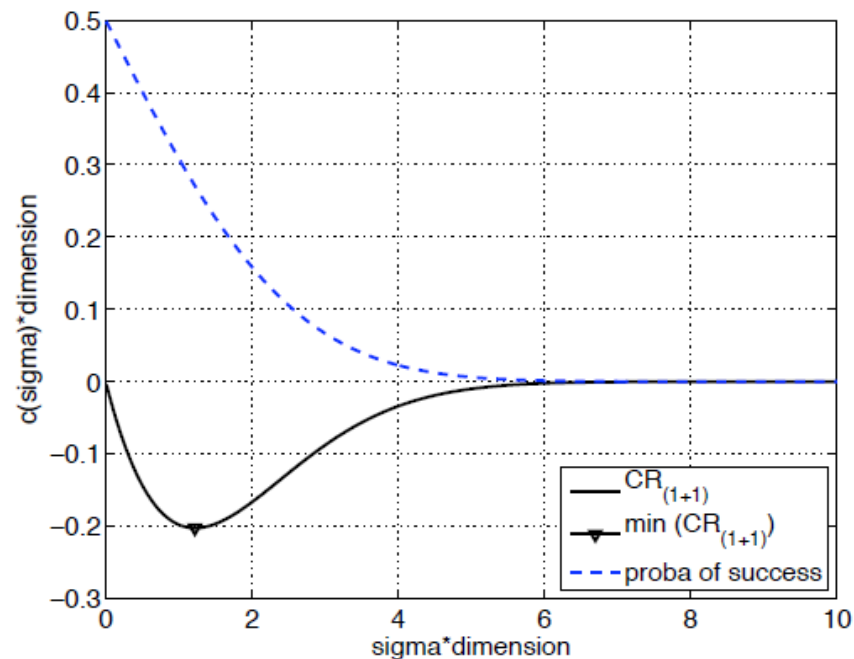
from [Auger, p. 34]



# 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)

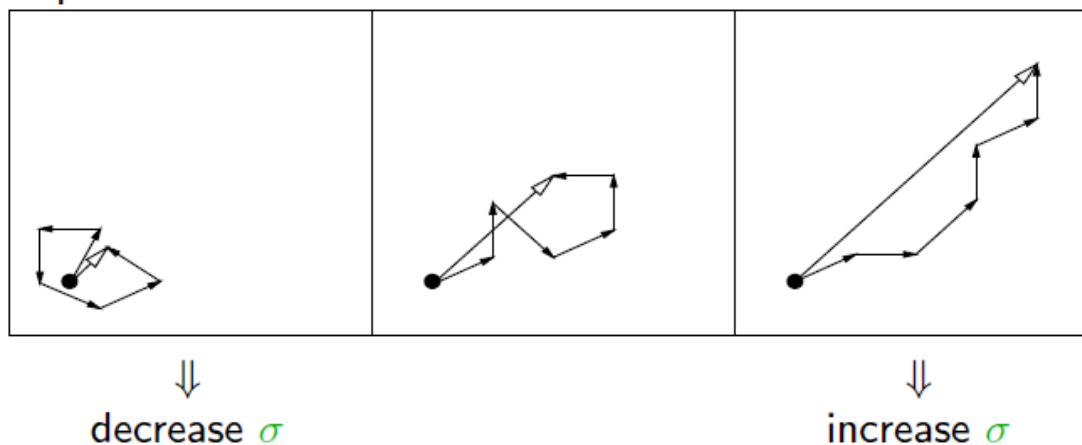
## Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

$$\begin{aligned}x_i &= m + \sigma y_i \\ m &\leftarrow m + \sigma y_w\end{aligned}$$

Measure the length of the *evolution path*

the pathway of the mean vector  $m$  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  $c_\sigma \approx 4/n$ ,  $d_\sigma \approx 1$ .

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

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu w}}_{\text{accounts for } w_i} \mathbf{y}_w$$

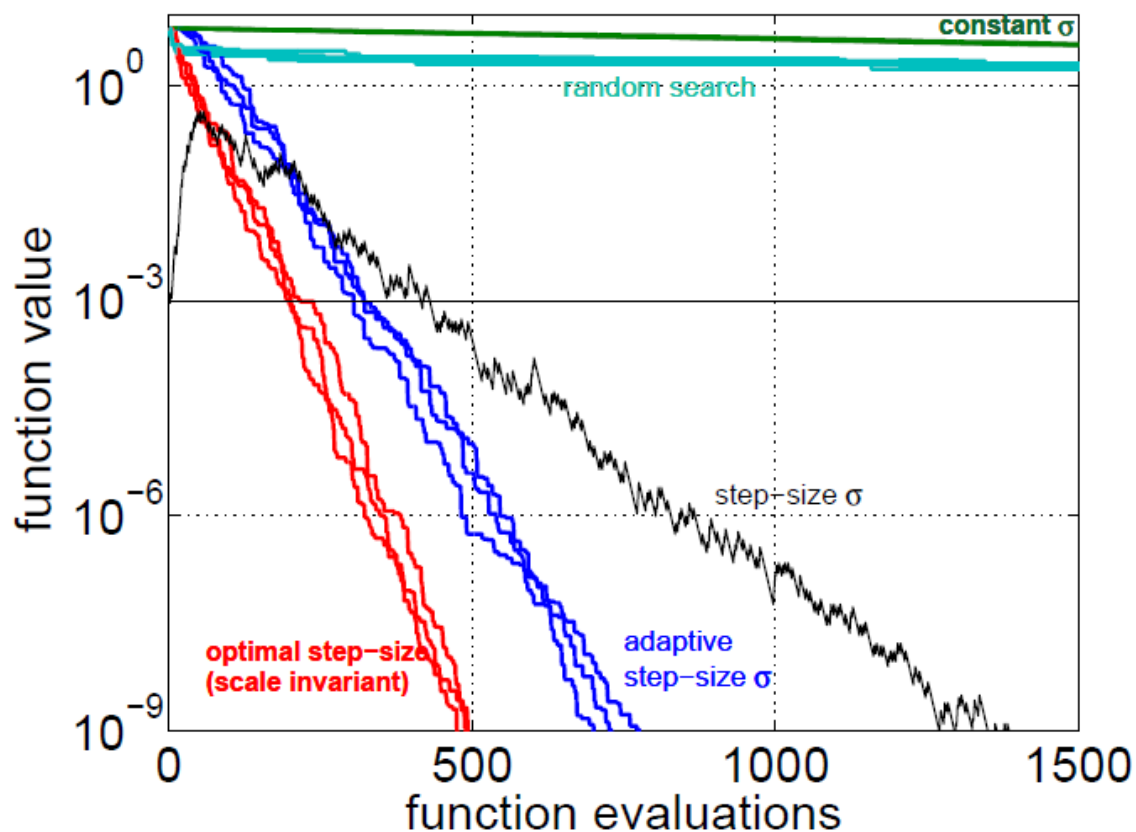
$$\sigma \leftarrow \sigma \times \underbrace{\exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right)}_{>1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}} \quad \text{update step-size}$$

from [Auger, p. 37]

# Cumulative Step-Size Adaptation (CSA)

## Step-size adaptation

What is achieved



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

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

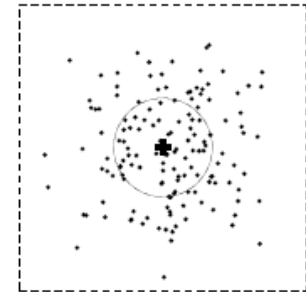
## Evolution Strategies

Recalling

New search points are sampled normally distributed

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

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
- ▶ 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  $\mathbf{C}$ .

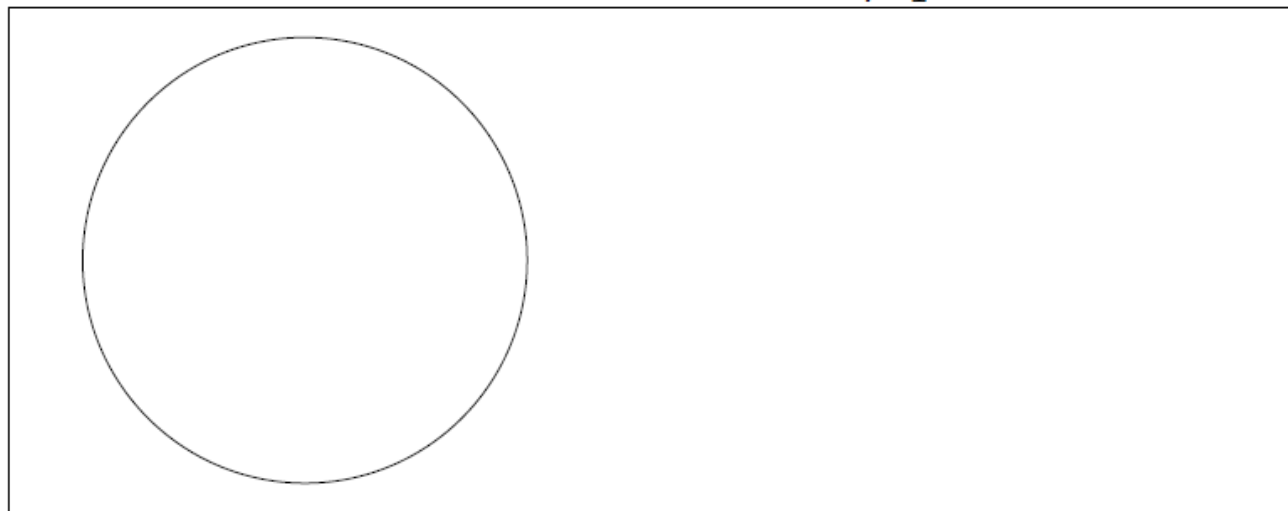
from [Auger, p. 40]

# Rank-One Update of Covariance Matrix

## Covariance Matrix Adaptation

### Rank-One Update

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



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

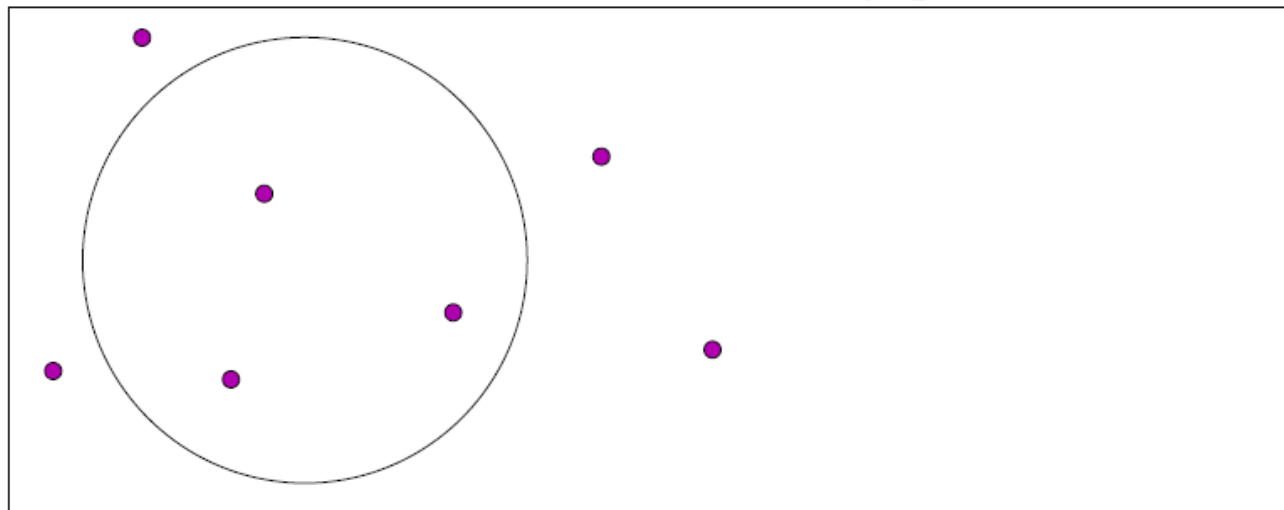
from [Auger, p. 41]

# Rank-One Update of Covariance Matrix

## Covariance Matrix Adaptation

### Rank-One Update

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



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

from [Auger, p. 41]

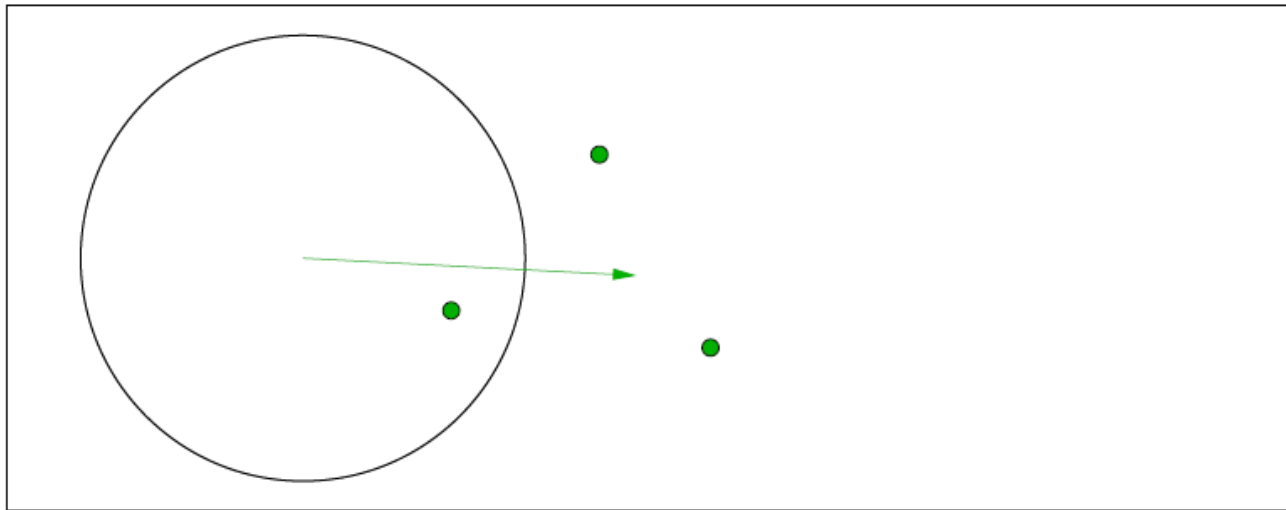


# Rank-One Update of Covariance Matrix

## Covariance Matrix Adaptation

### Rank-One Update

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



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

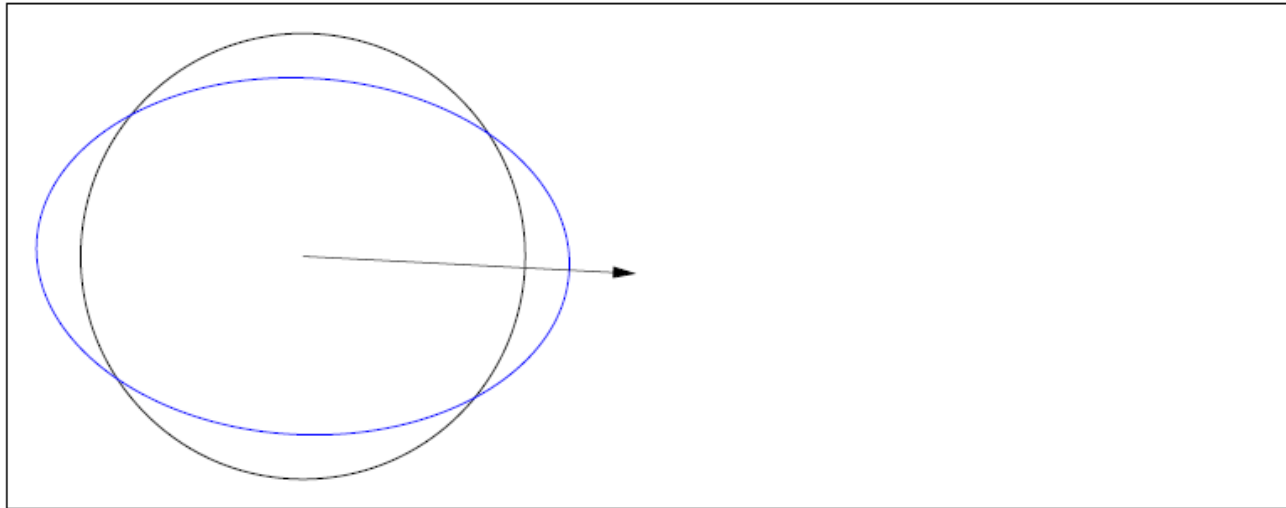
from [Auger, p. 41]

# Rank-One Update of Covariance Matrix

## Covariance Matrix Adaptation

### Rank-One Update

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



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

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

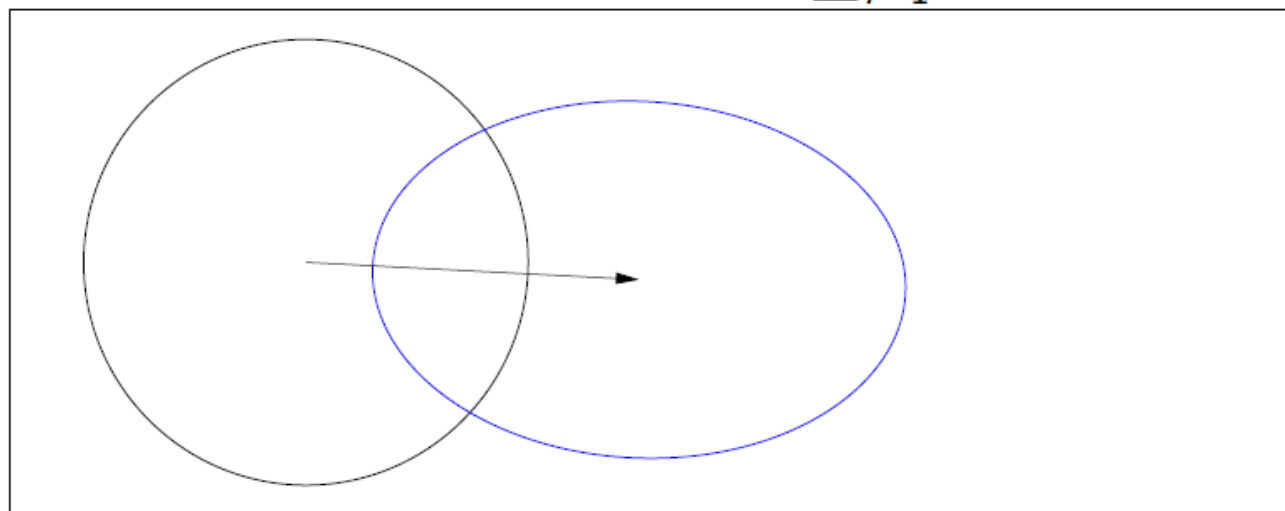
from [Auger, p. 41]

# 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  $\sigma$ )

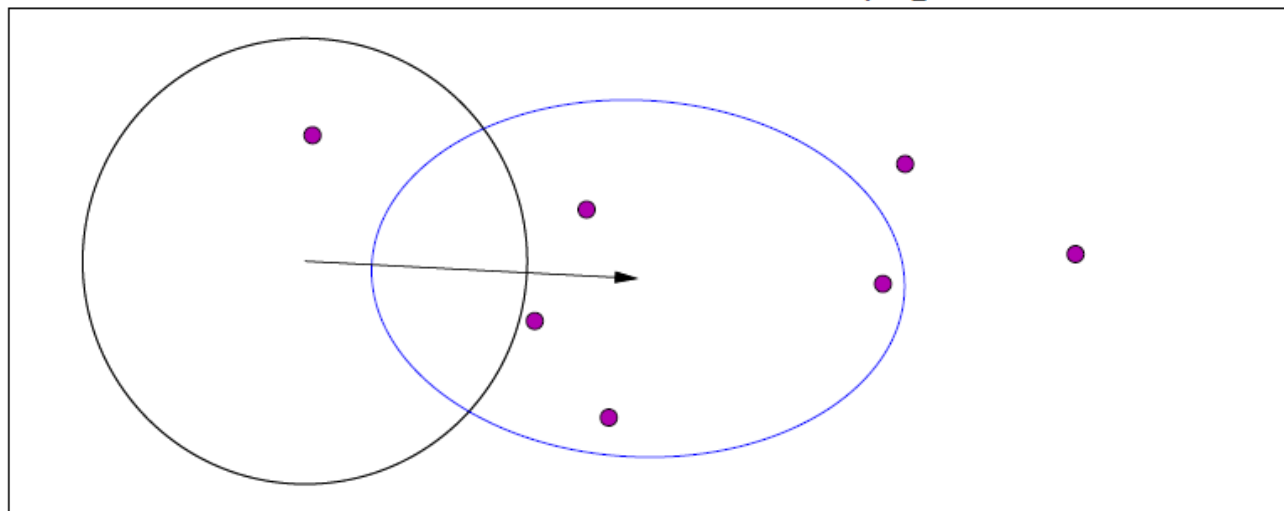
from [Auger, p. 41]

# 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  $\sigma$ )

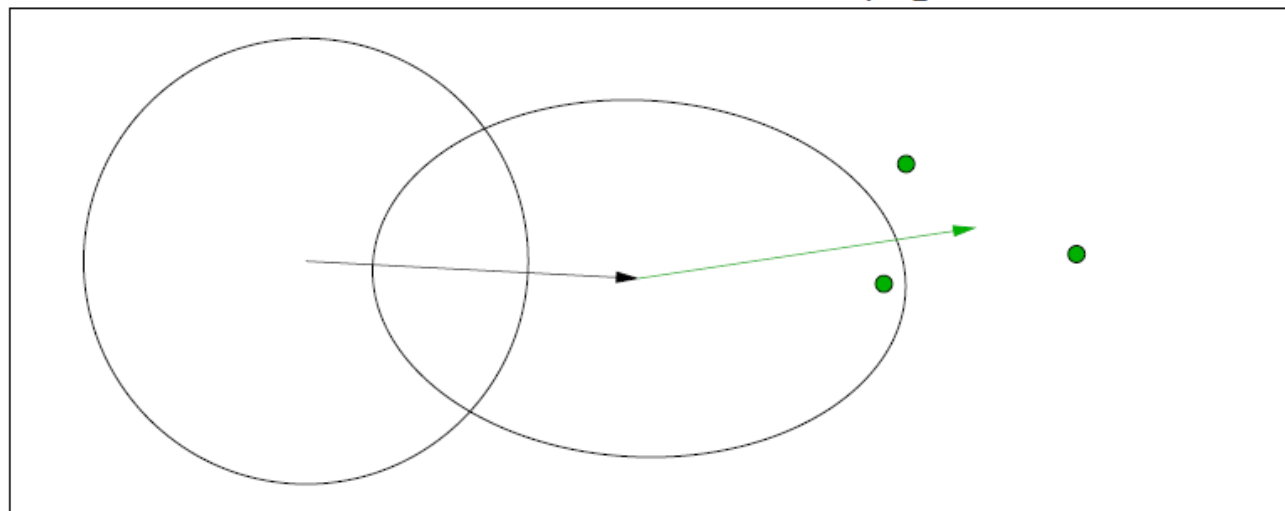
from [Auger, p. 41]

# 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  $\mathbf{m}$

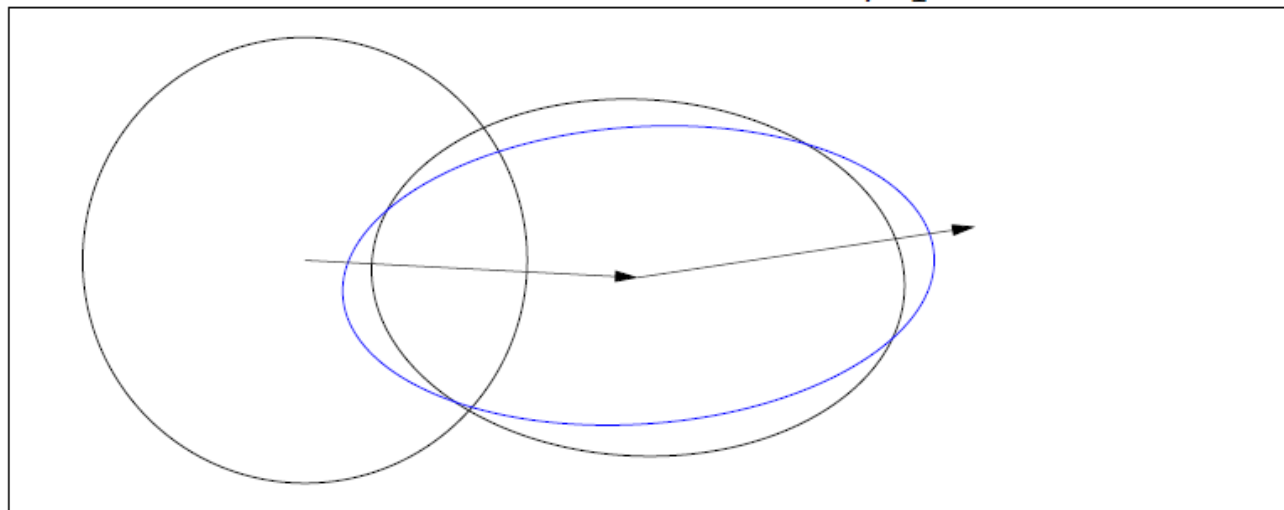
from [Auger, p. 41]

# Rank-One Update of Covariance Matrix

## Covariance Matrix Adaptation

### Rank-One Update

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



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

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

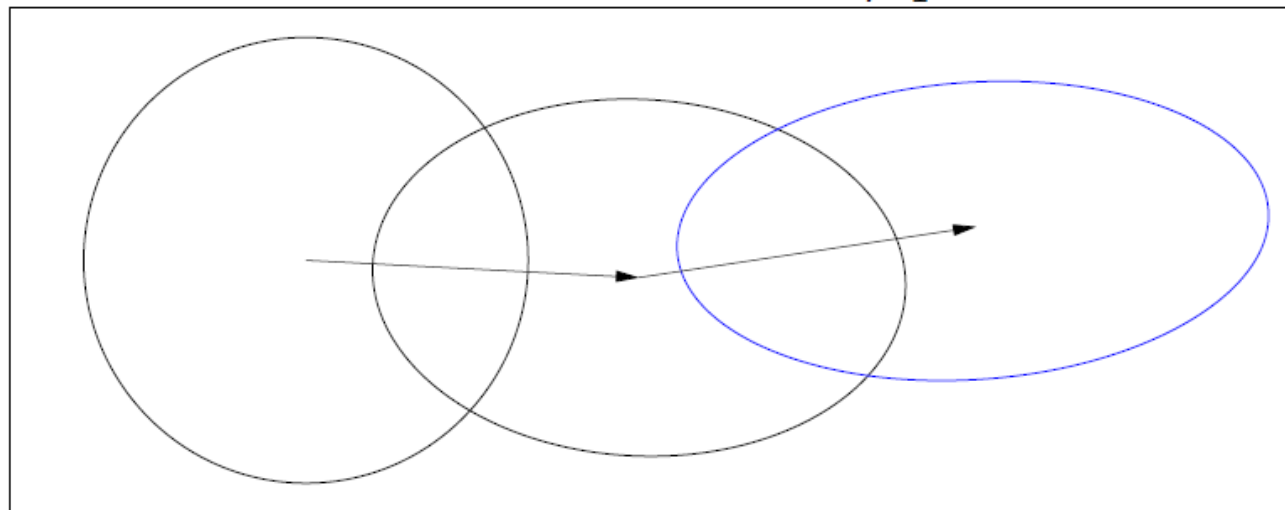
from [Auger, p. 41]

# 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,

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

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

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

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

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

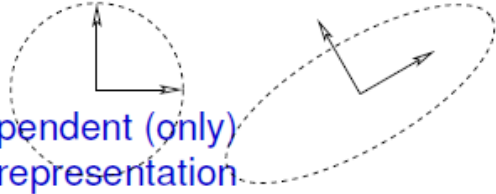
from [Auger, p. 42]



# 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

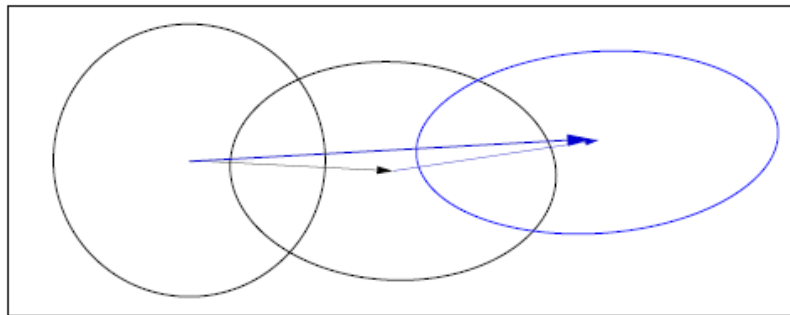
# Evolution Path

## 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  $m$ .



An exponentially weighted sum of steps  $y_w$  is used

$$p_c \propto \sum_{i=0}^g \underbrace{(1 - c_c)^{g-i}}_{\text{exponentially fading weights}} y_w^{(i)}$$

The recursive construction of the evolution path (cumulation):

$$p_c \leftarrow \underbrace{(1 - c_c)}_{\text{decay factor}} p_c + \underbrace{\sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w}}_{\text{normalization factor}} \underbrace{y_w}_{\text{input} = \frac{m - m_{\text{old}}}{\sigma}}$$

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

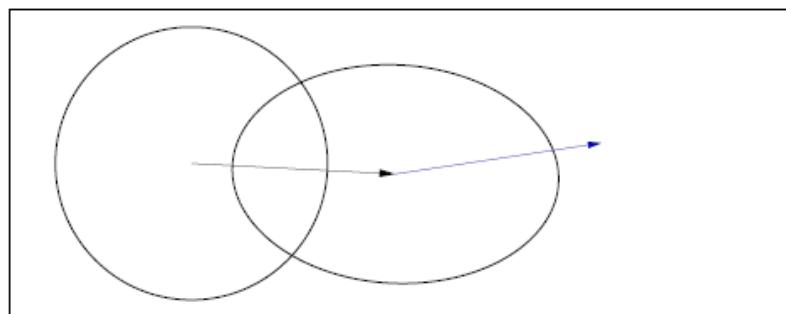
from [Auger, p. 44]

# Utilizing the Evolution Path

## Cumulation

### Utilizing the Evolution Path

We used  $\mathbf{y}_w \mathbf{y}_w^T$  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.



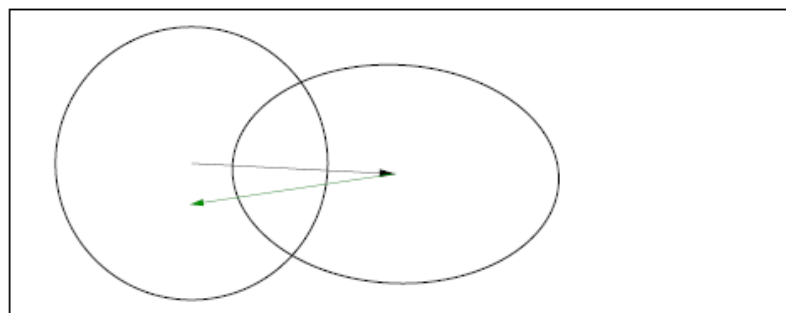
from [Auger, p. 45]

# 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.



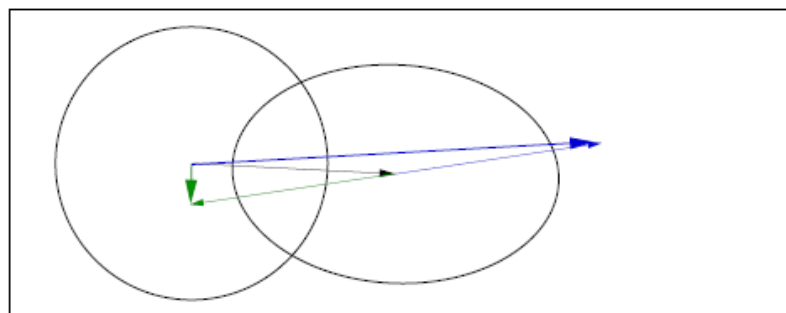
from [Auger, p. 45

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

from [Auger, p. 45]

# Rank- $\mu$ Update

## Rank- $\mu$ 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- $\mu$  update extends the update rule for **large population sizes**  $\lambda$  using  $\mu > 1$  vectors to update  $\mathbf{C}$  at each generation step.

from [Auger, p. 47]

# Rank- $\mu$ Update

## Rank- $\mu$ 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- $\mu$  update extends the update rule for **large population sizes**  $\lambda$  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  $\mu$  steps and has rank  $\min(\mu, n)$  with probability one.

from [Auger, p. 47]

# Rank- $\mu$ Update

## Rank- $\mu$ 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- $\mu$  update extends the update rule for **large population sizes**  $\lambda$  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  $\mu$  steps and has rank  $\min(\mu, n)$  with probability one.

The rank- $\mu$  update then reads

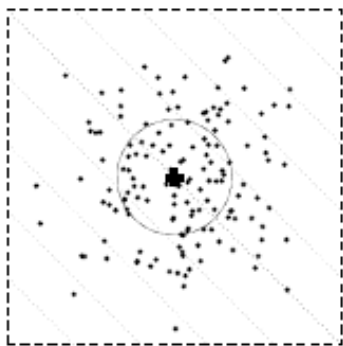
$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_{\mu}$$

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

from [Auger, p. 47]



# Illustration of Rank- $\mu$ Update



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

sampling of

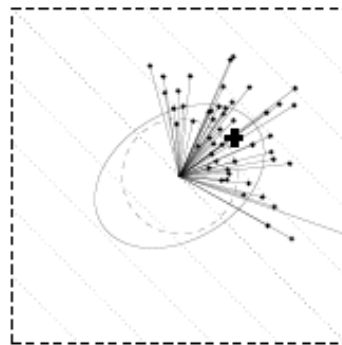
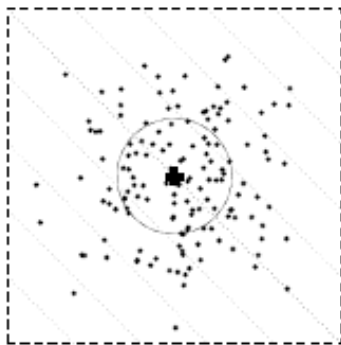
$\lambda = 150$  solutions

where  $\mathbf{C} = \mathbf{I}$  and

$$\sigma = 1$$

from [Auger, p. 48]

# Illustration of Rank- $\mu$ Update



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

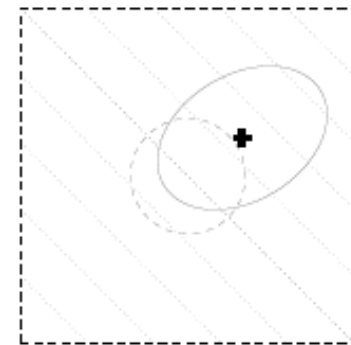
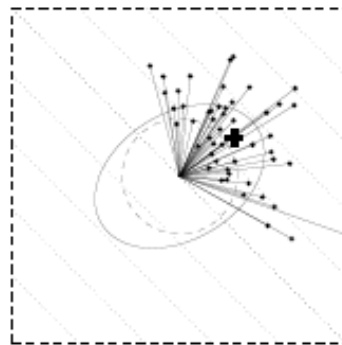
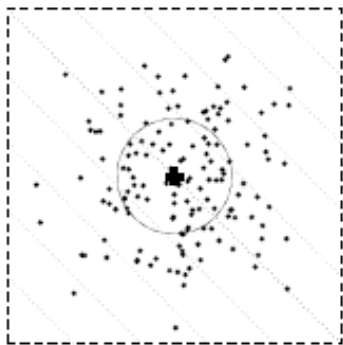
$$\mathbf{C} \leftarrow \frac{1}{(\mu - 1) \times \mathbf{C} + 1} \times \mathbf{C}_\mu$$

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  
 $\mathbf{C}_{\text{COV}} = 1$

from [Auger, p. 48]

# Illustration of Rank- $\mu$ Update



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

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

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

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

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

new distribution

from [Auger, p. 48]

# Rank- $\mu$ Update: Summary

## The rank- $\mu$ 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

## The CMA-ES

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

**Initialize:**  $C = I$ , and  $p_c = \mathbf{0}$ ,

**Set:**  $c_c \approx 4/n$ ,  $c_\sigma \approx 4/n$ ,  $c_1 \approx 1/n$

and  $w_{i=1\dots\lambda}$  such that  $\sum w_i = 1$

### Promised:

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

$\frac{w_i}{n}$ ,

**While not terminate**

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

$$m \leftarrow \sum_{i=1}^{\mu} w_i x_{i:\lambda} = m + \sigma y_w \quad \text{where } y_w = \sum_{i=1}^{\mu} w_i y_{i:\lambda} \quad \text{update mean}$$

$$p_c \leftarrow (1 - c_c) p_c + \mathbb{1}_{\{\|p_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} y_w \quad \text{cumulation for } C$$

$$p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} C^{-\frac{1}{2}} y_w \quad \text{cumulation for } \sigma$$

$$C \leftarrow (1 - c_1 - c_\mu) C + c_1 p_c p_c^T + c_\mu \sum_{i=1}^{\mu} w_i y_{i:\lambda} y_{i:\lambda}^T \quad \text{update } C$$

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

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

## The CMA-ES

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

**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**

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

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

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

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

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

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

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

## Strategy Internal Parameters

- related to selection and recombination
  - ▶  $\lambda$ , offspring number, new solutions sampled, population size
  - ▶  $\mu$ , parent number, solutions involved in updates of  $m$ ,  $\mathbf{C}$ , and  $\sigma$
  - ▶  $w_{i=1,\dots,\mu}$ , recombination weights
- related to  $\mathbf{C}$ -update
  - ▶  $c_c$ , decay rate for the evolution path
  - ▶  $c_1$ , learning rate for rank-one update of  $\mathbf{C}$
  - ▶  $c_\mu$ , learning rate for rank- $\mu$  update of  $\mathbf{C}$
- related to  $\sigma$ -update
  - ▶  $c_\sigma$ , decay rate of the evolution path
  - ▶  $d_\sigma$ , damping for  $\sigma$ -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  $\lambda$  (and the initial  $\sigma$ ) might be reasonably varied in a wide range, *depending on the objective function*

Useful: restarts with increasing population size (IPOP)

# Experimental Considerations



## 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

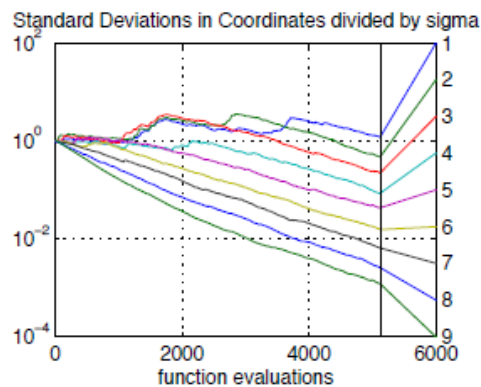
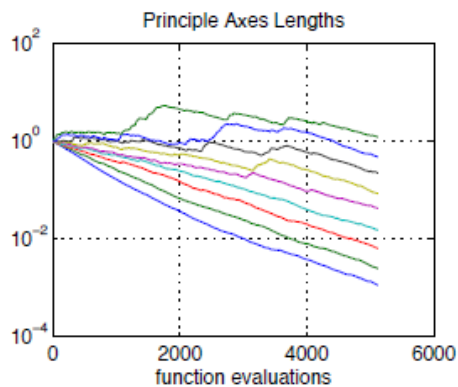
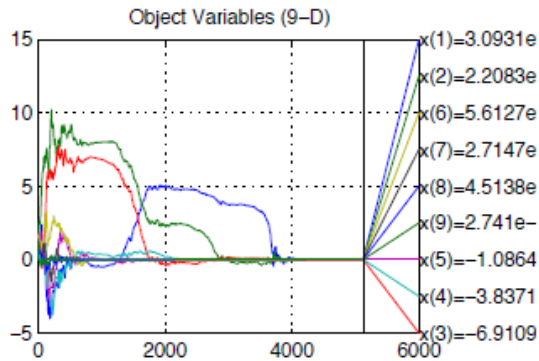
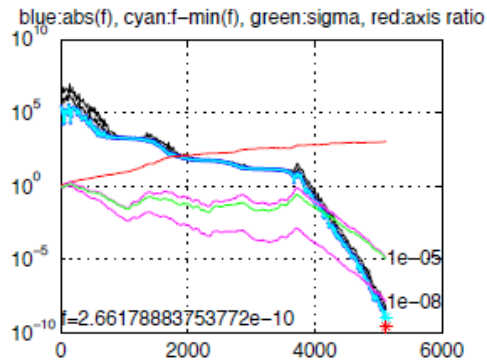
$$\mathbf{C} \propto \mathbf{H}^{-1}$$

in a stochastic sense

# Experimentum Crucis with CMA-ES

## Experimentum Crucis (1)

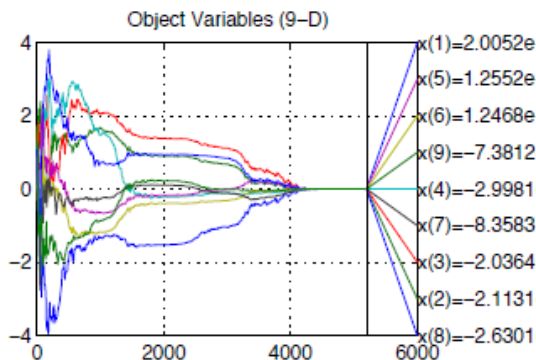
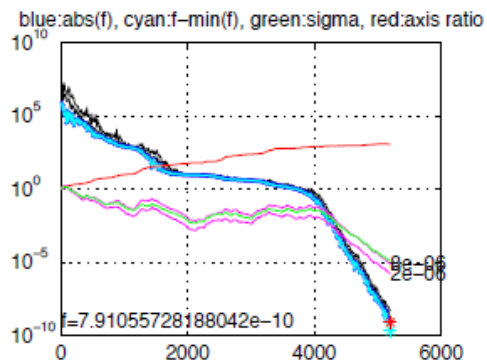
$f$  convex quadratic, separable



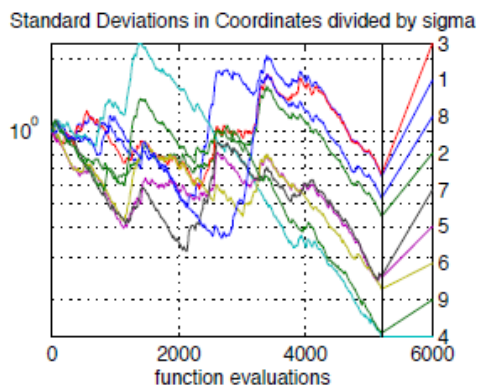
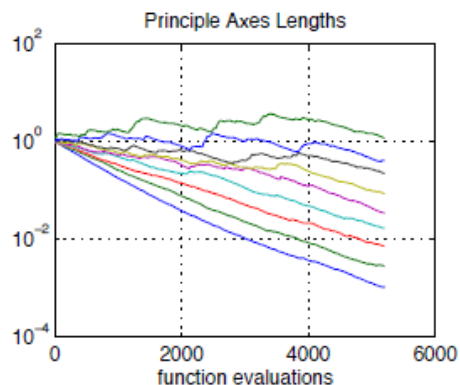
$$f(\mathbf{x}) = \sum_{i=1}^n 10^{\alpha \frac{i-1}{n-1}} x_i^2, \alpha = 6$$

## Experimentum Crucis (2)

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



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



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

from [Hansen, p. 93]

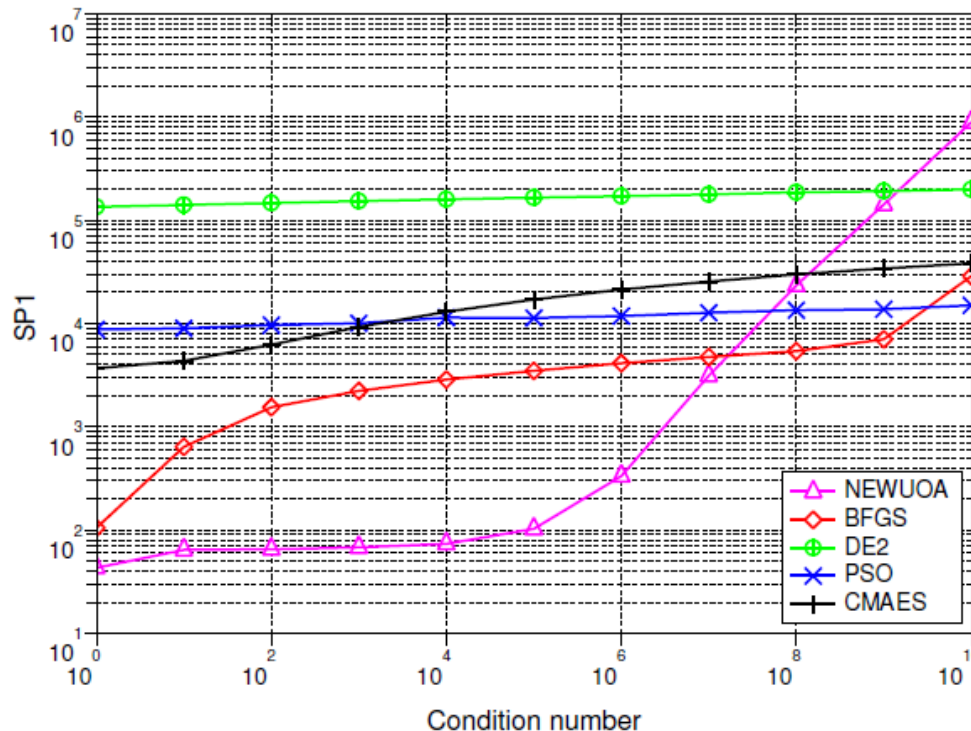
# Influence of Condition Number + Invariance

## Comparing Experiments

### Comparison to BFGS, NEWUOA, PSO and DE

$f$  convex quadratic, separable with varying condition number  $\alpha$

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)

$f(x) = g(x^T H x)$  with

$H$  diagonal

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

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

SP1 = average number of objective function evaluations<sup>14</sup> to reach the target function value of  $g^{-1}(10^{-9})$

<sup>14</sup> Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

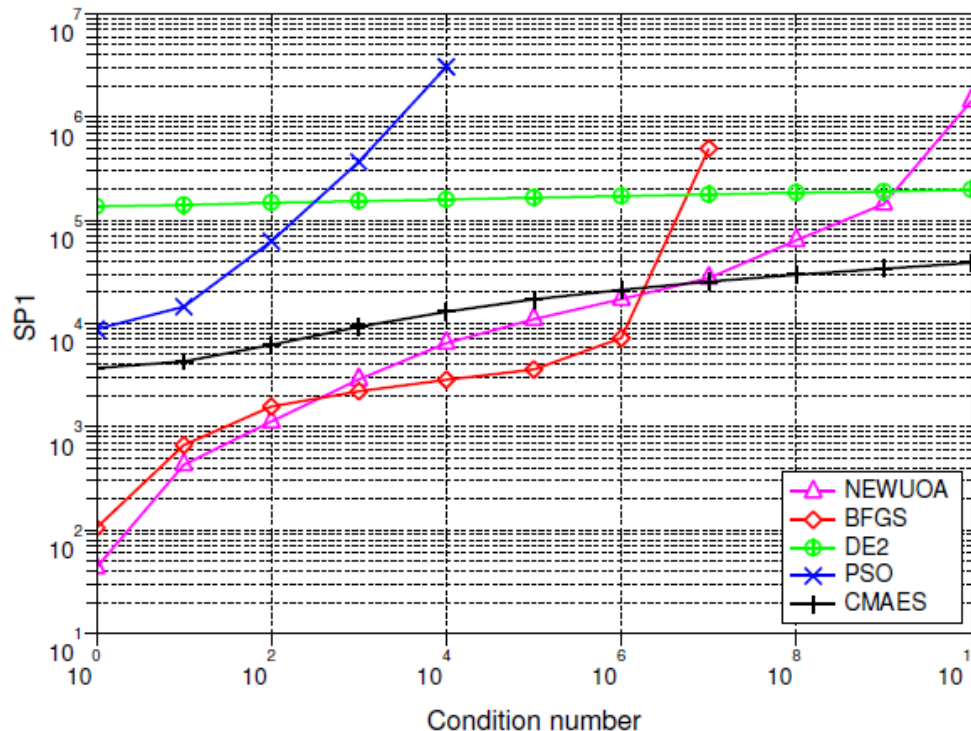
# Influence of Condition Number + Invariance

## Comparing Experiments

### Comparison to BFGS, NEWUOA, PSO and DE

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

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)

$f(x) = g(x^T H x)$  with

$H$  full

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

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

SP1 = average number of objective function evaluations<sup>15</sup> to reach the target function value of  $g^{-1}(10^{-9})$

<sup>15</sup> Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

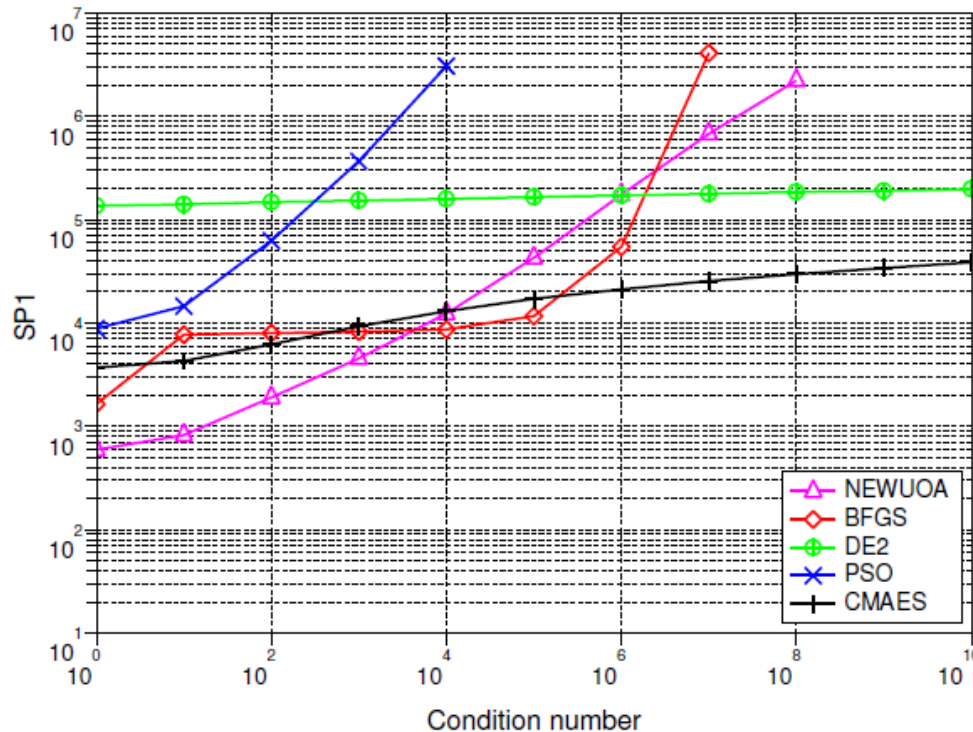
# Influence of Condition Number + Invariance

## Comparing Experiments

### Comparison to BFGS, NEWUOA, PSO and DE

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

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)

$f(x) = g(x^T H x)$  with

$H$  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 evaluations<sup>16</sup> to reach the target function value of  $g^{-1}(10^{-9})$

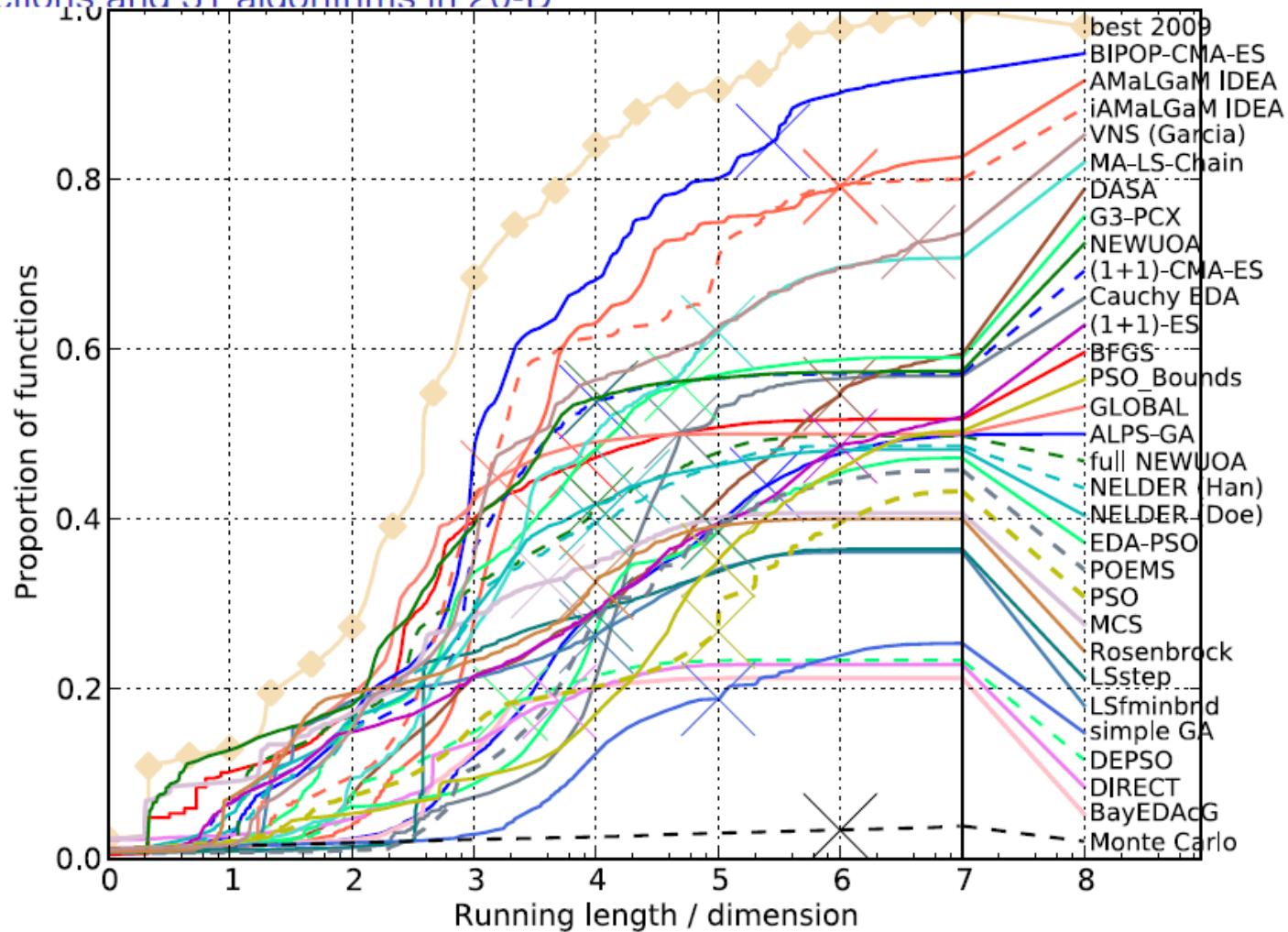
<sup>16</sup>Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

# 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