

Advanced Optimization

Lecture 6: Stochastic Algorithms II (SGD & CMA-ES)

November 17, 2021

CentraleSupélec / ESSEC Business School

dimo.brockhoff@inria.fr



Dimo Brockhoff
Inria Saclay – Ile-de-France



INSTITUT
POLYTECHNIQUE
DE PARIS



Course Overview

		Topic
Wed, 13.10.2021	PM	Introduction, examples of problems, problem types
Wed, 20.10.2021	PM	Continuous (unconstrained) optimization: convexity, gradients, Hessian, ... [technical test Evalmee]
Wed, 27.10.2021	PM	Continuous optimization II: [1 st mini-exam] Constrained optimization: Lagrangian, optimality conditions
Wed, 03.11.2021	PM	gradient descent, Newton direction, quasi-Newton (BFGS) Linear programming: duality, maxflow/mincut, simplex algo
Wed, 10.11.2021	PM	derivative-free algorithms: Nelder-Mead and CMA-ES
Wed, 17.11.2021	PM	CMA-ES, Part II, Stochastic Gradient Descent, Bayesian optimization
Wed, 24.11.2021	PM	Benchmarking solvers: runtime distributions, performance profiles [2nd mini-exam]
Tue, 30.11.2021	23:59	Deadline open source project (PDF sent by email)
Wed, 01.12.2021	PM	Discrete optimization: branch and bound, branch and cut, k-means clustering
Wed, 15.12.2021	PM	Exam

Details on Continuous Optimization Lectures

Introduction to Continuous Optimization

- examples and typical difficulties in optimization

Mathematical Tools to Characterize Optima

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

Gradient-based Algorithms

- gradient descent
- quasi-Newton method (BFGS) and invariances

Linear programming, duality

Learning in Optimization / Optimization in Machine Learning

- Stochastic gradient descent (SGD) + Adam
- CMA-ES (adaptive algorithms / Information Geometry)
- Other derivative-free algorithms: Nelder-Mead, Bayesian opt.

back to

Covariance Matrix Adaptation Evolution Strategy (CMA-ES)

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

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} \mathbf{y}_i \mathbf{y}_i^T$ 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.

Copyright Notice

- As in the last lecture, slides are taken partly from <http://www.cmap.polytechnique.fr/~nikolaus.hansen/copenhagen-cma-es.pdf> (copyright by Nikolaus Hansen) and from <http://www.cmap.polytechnique.fr/~dimo.brockhoff/optimizationSaclay/2015/slides/20151106-continuousoptIV.pdf> (copyright by Anne Auger)
- Those are referred to as [Hansen, p. X] and [Auger, p. Y] respectively.

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 $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$
- 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

Step-Size Adaptation

Recap CMA-ES: What We Have So Far

Step-Size Control

Evolution Strategies

Recalling

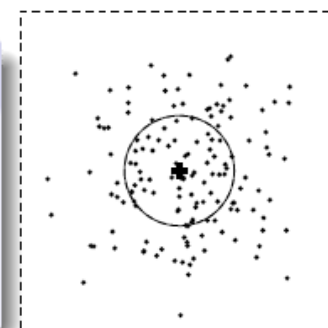
New search points are sampled normally distributed

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

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

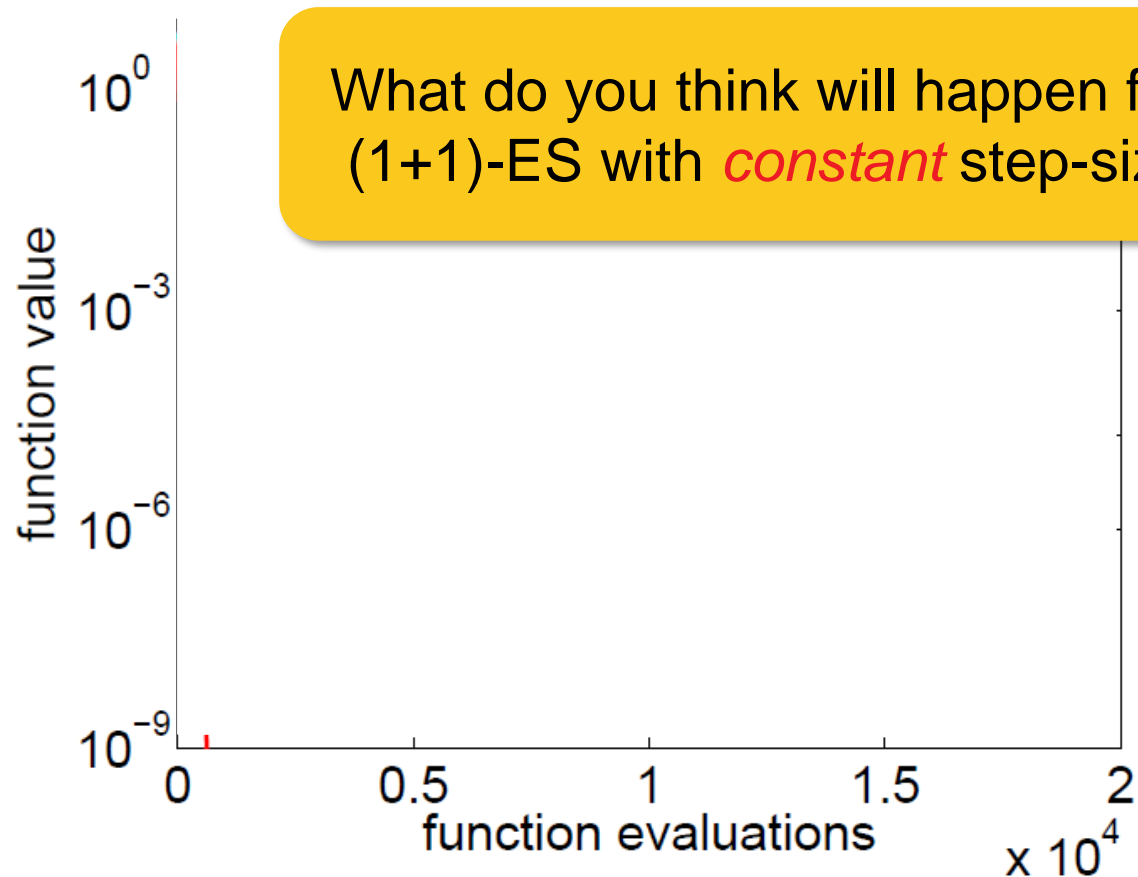


The remaining question is how to update σ and \mathbf{C} .

from [Hansen, p. 45]

Why At All Step-Size Adaptation?

Why Step-Size Control?



What do you think will happen for a (1+1)-ES with *constant* step-size?

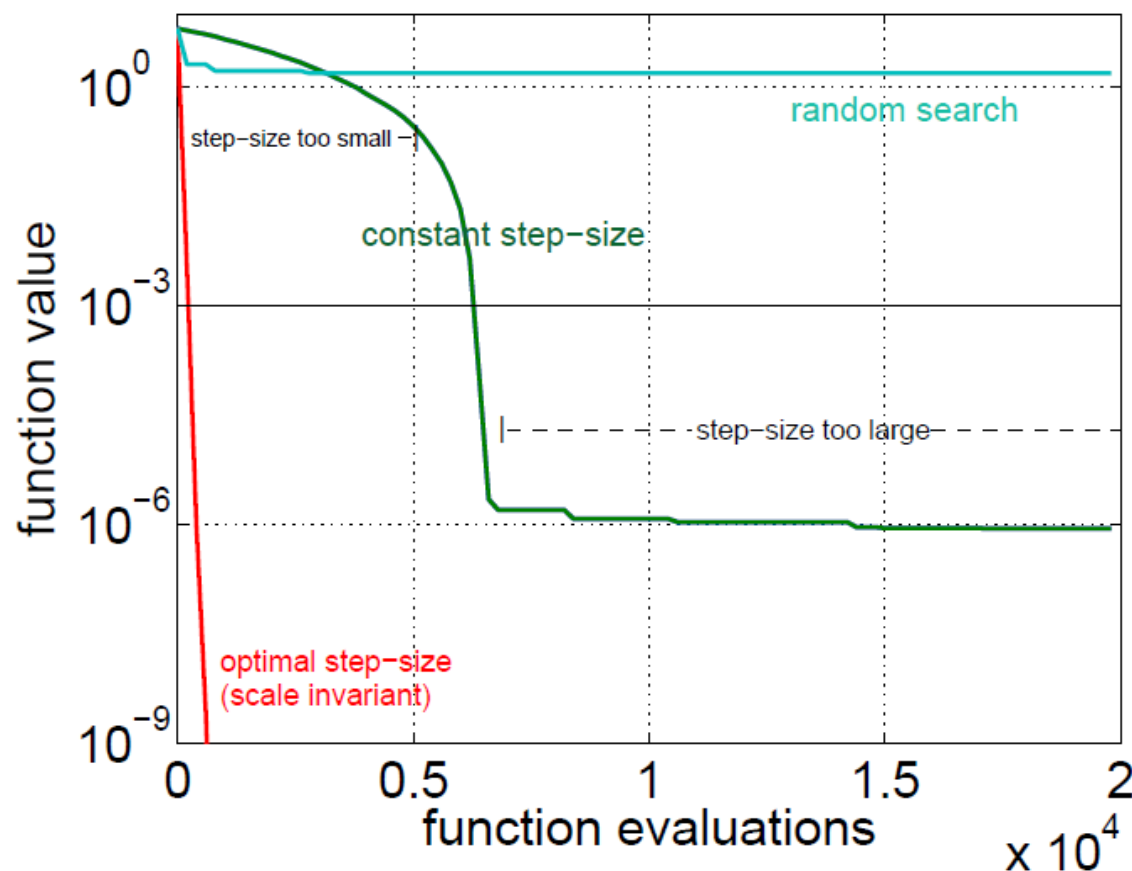
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

from [Auger, p. 22]

Why Step-Size Adaptation?

Why Step-Size Control?



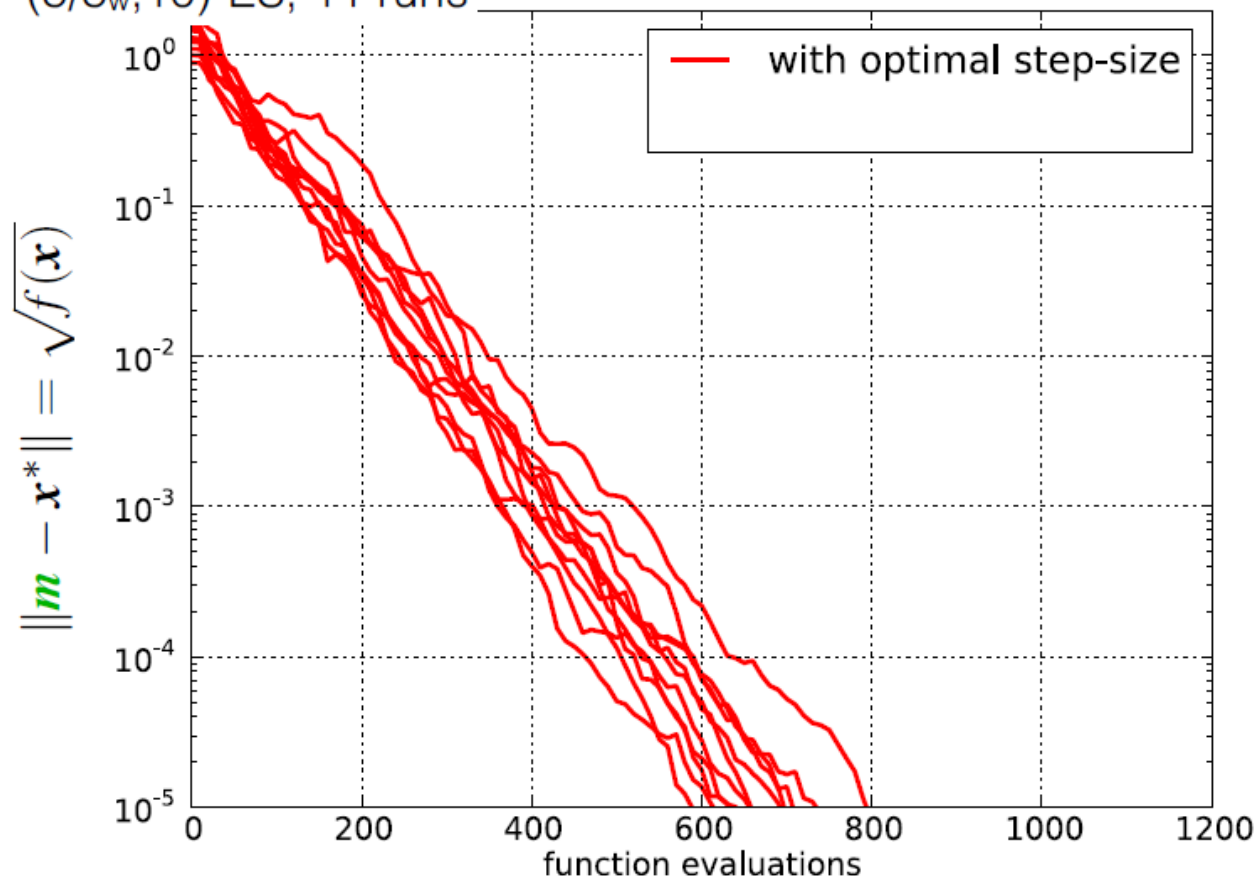
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

from [Auger, p. 22]

Why Step-Size Control?

(5/5_w,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 σ

from [Hansen, p. 47]

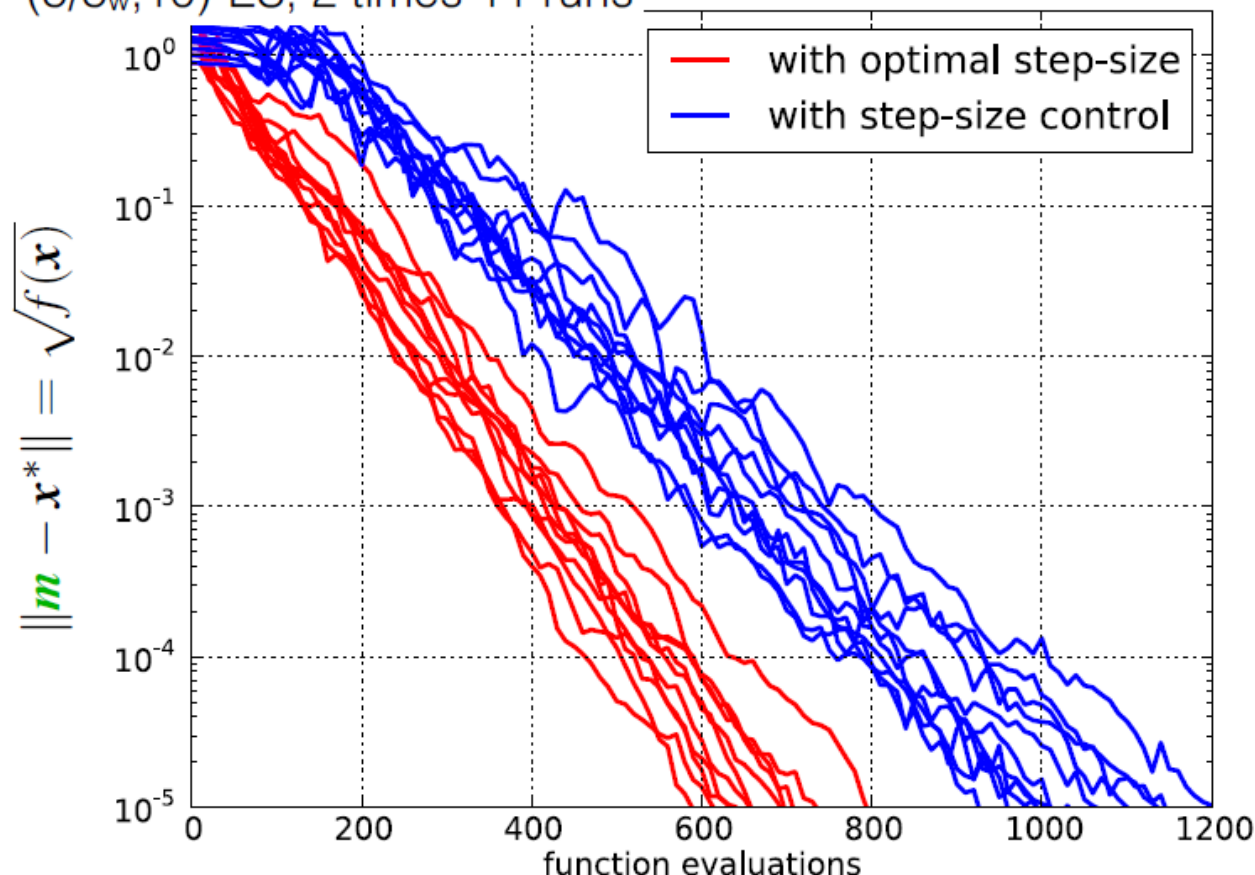
Optimal Step-Size vs. Step-Size Control

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]

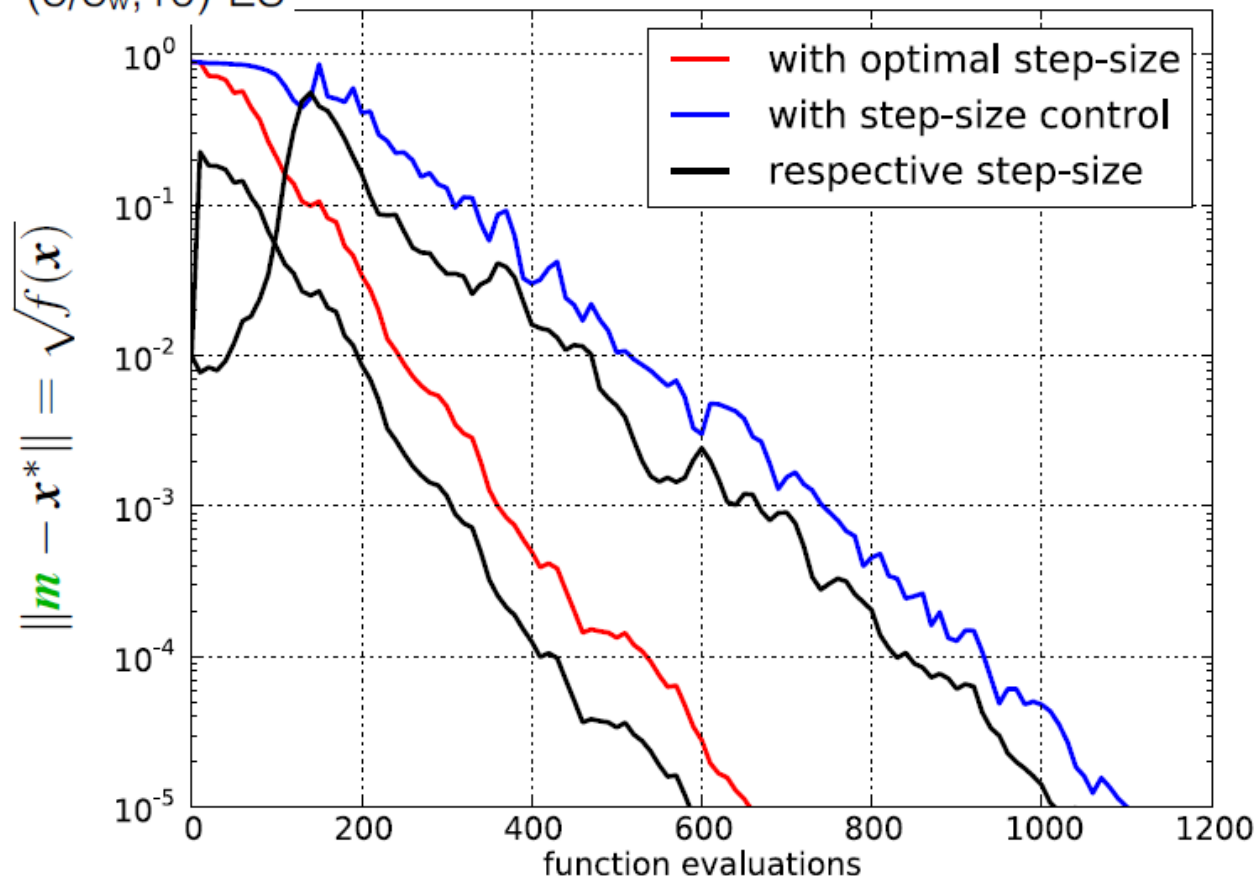
Optimal Step-Size vs. Step-Size Control

Step-Size Control

Why Step-Size Control

Why Step-Size Control?

(5/5_w, 10)-ES



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

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

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:

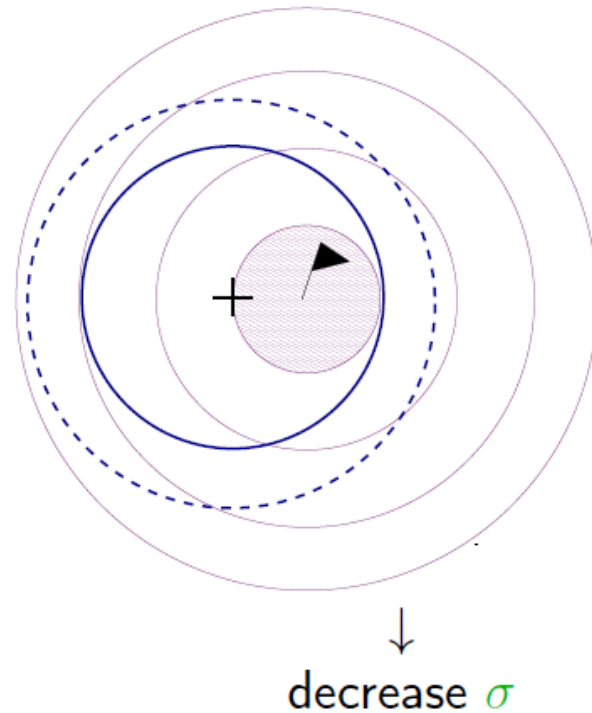
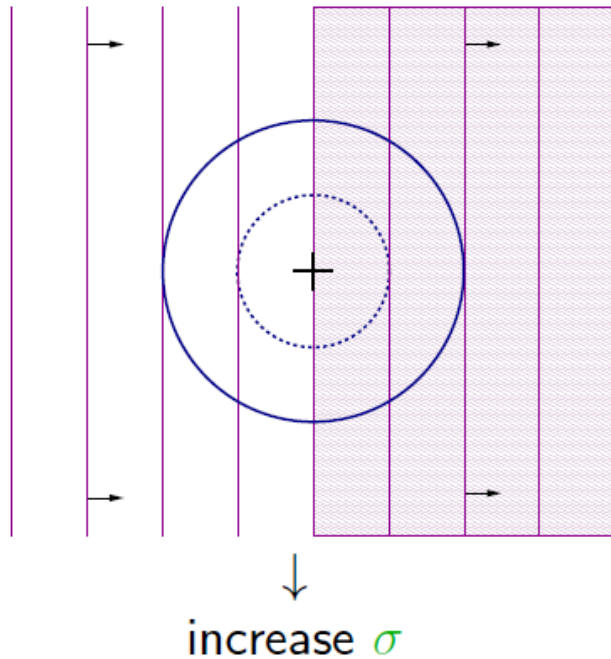
How to actually adapt the step-size during the optimization?

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

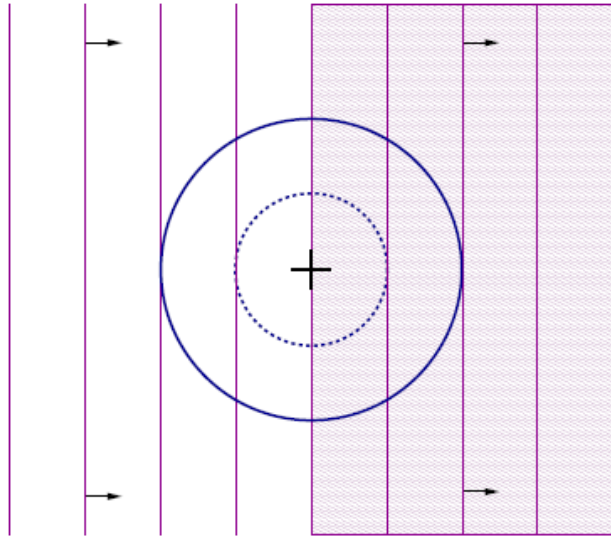
One-fifth success rule



from [Auger, p. 32]

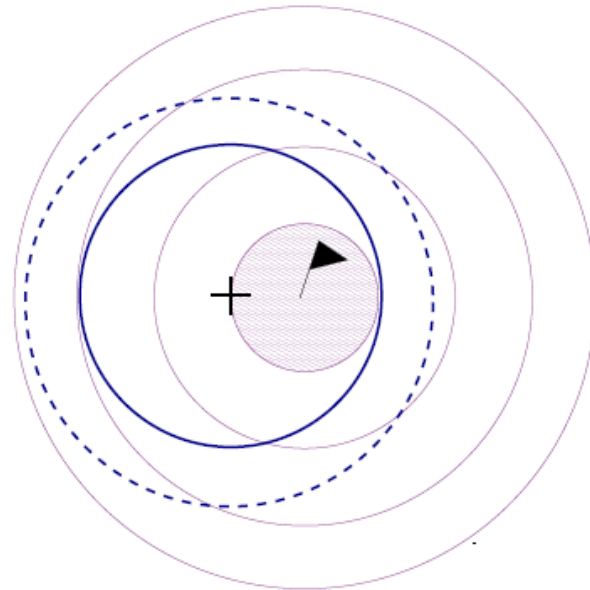
One-Fifth Success Rule

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

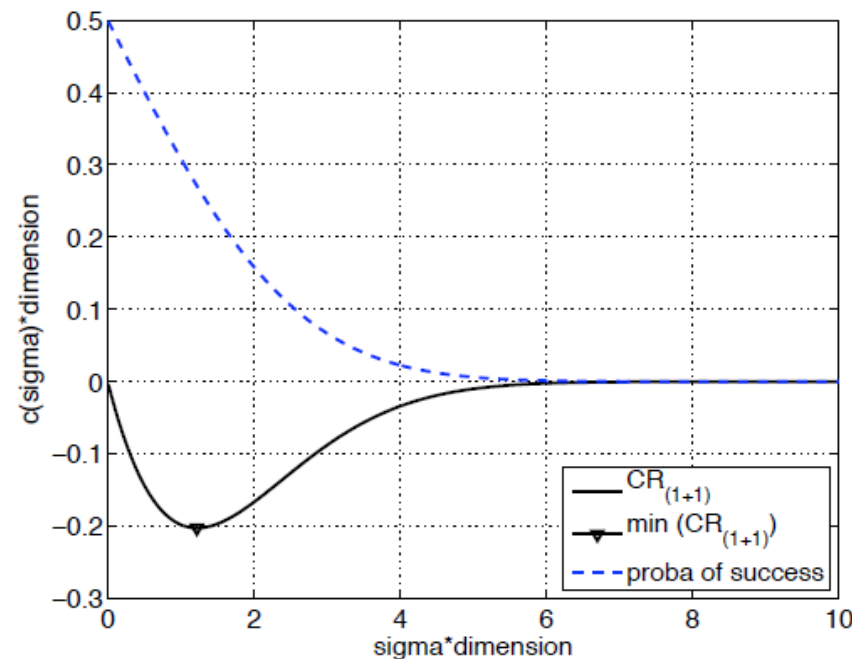
$$p_s = 0, \sigma \leftarrow \sigma / \exp(1/3)^{1/4}$$

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)

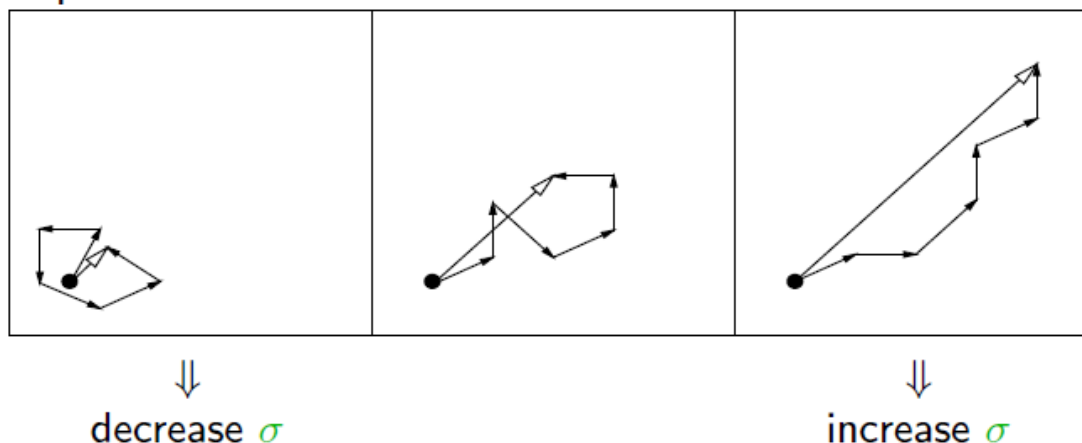
Path Length Control (CSA)

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 $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

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

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

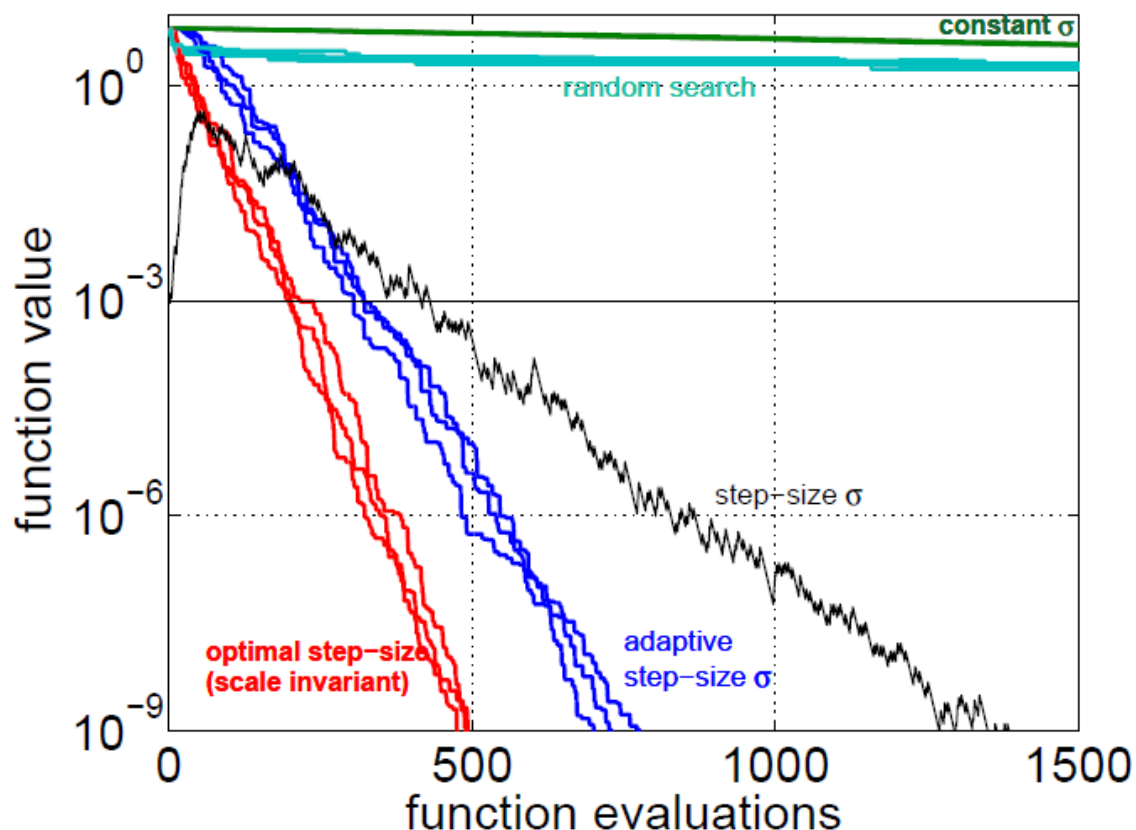
$$\sigma \leftarrow \sigma \times \underbrace{\exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)}_{>1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}} \quad \text{update step-size}$$

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

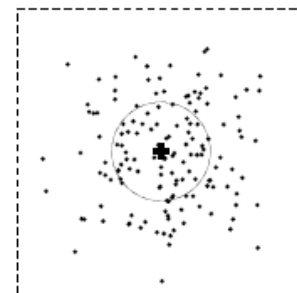
Evolution Strategies

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} .

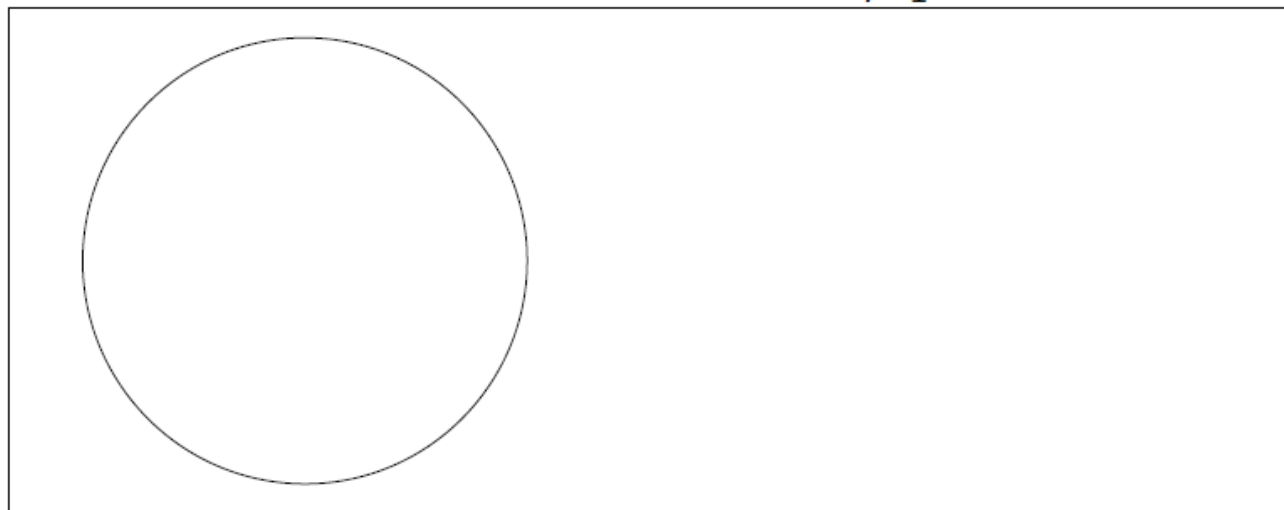
from [Auger, p. 40]

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

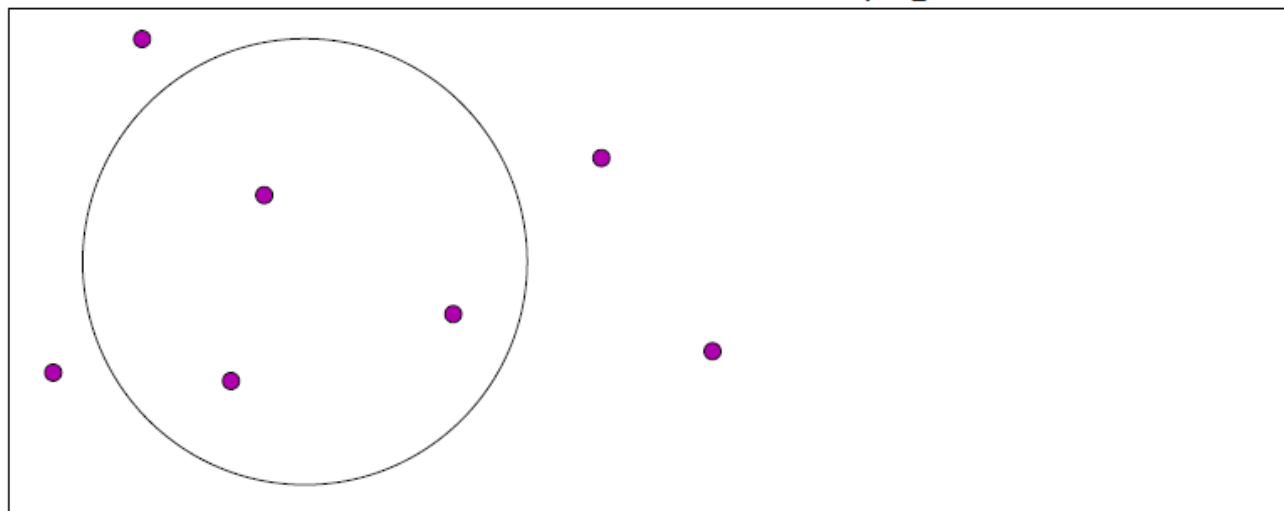
from [Auger, p. 41]

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

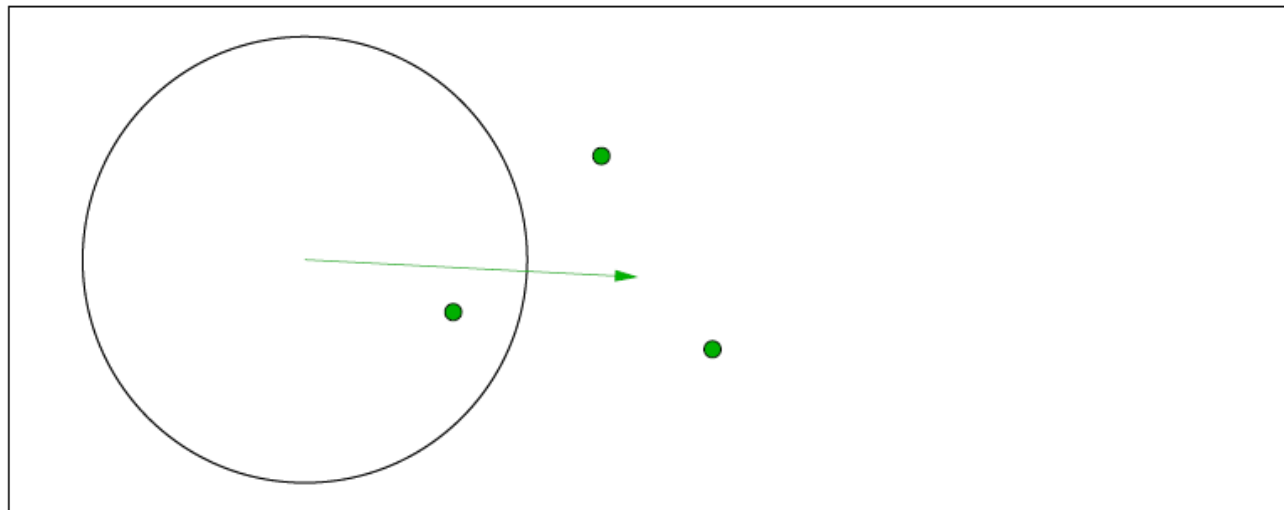
from [Auger, p. 41]

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 σ)

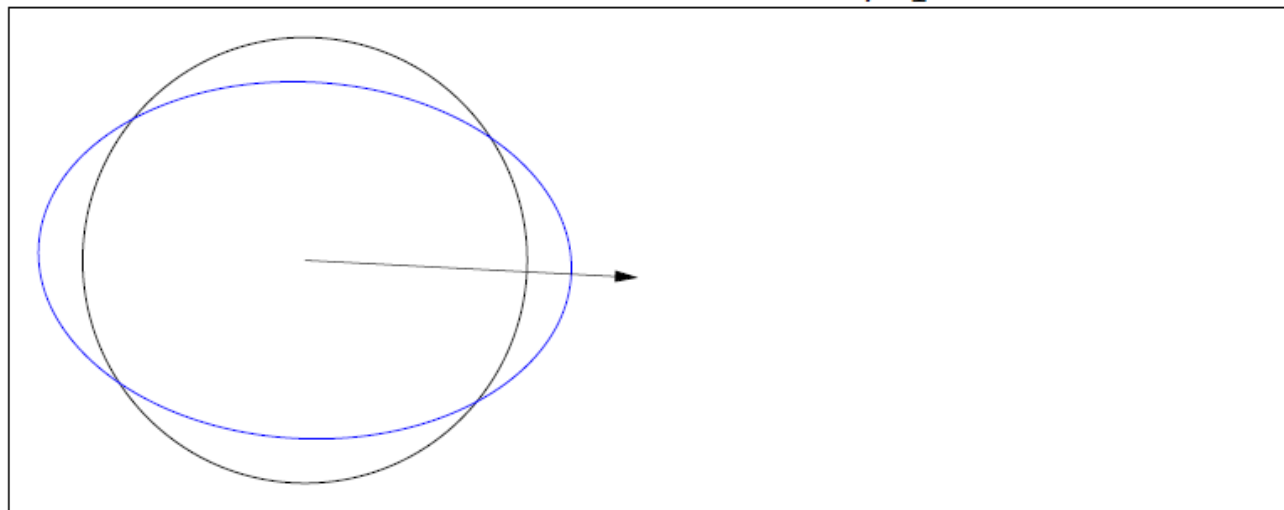
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



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

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

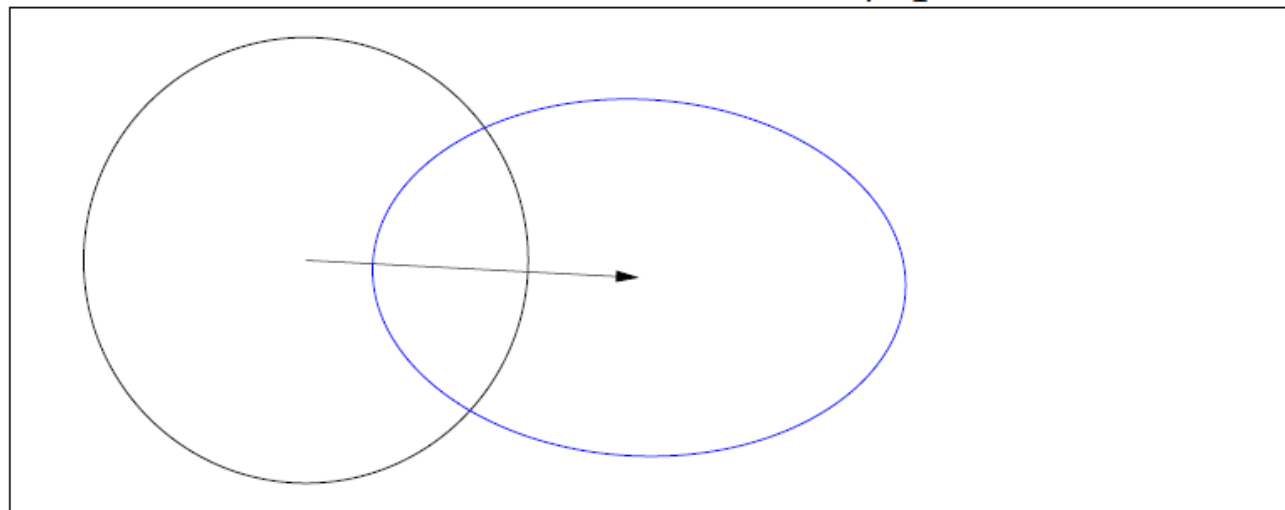
from [Auger, p. 41]

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)

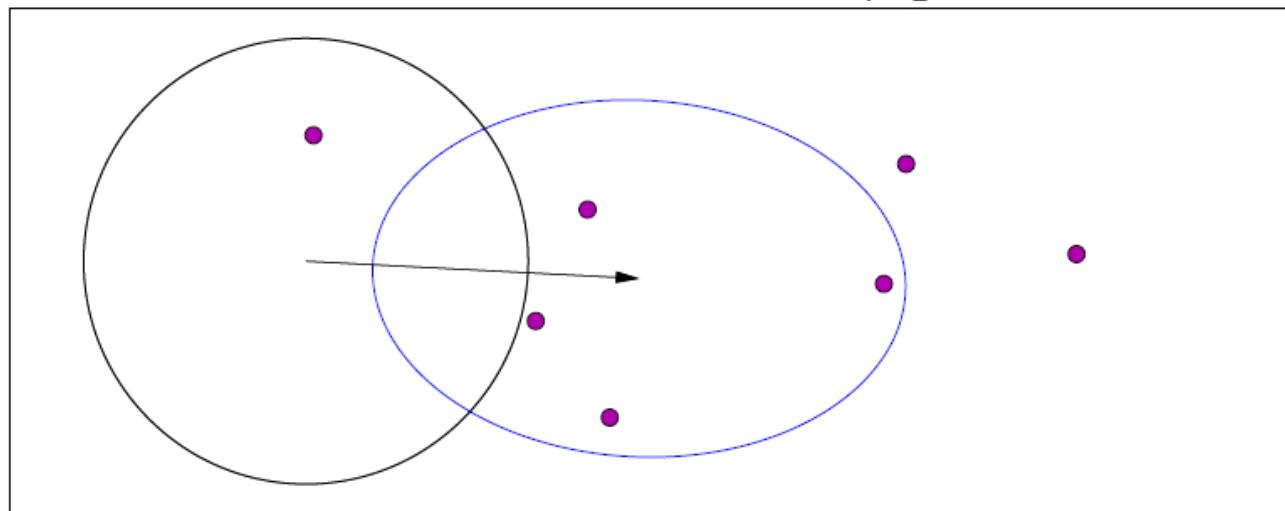
from [Auger, p. 41]

Rank-One Update of Covariance Matrix

Covariance Matrix Adaptation

Rank-One Update

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



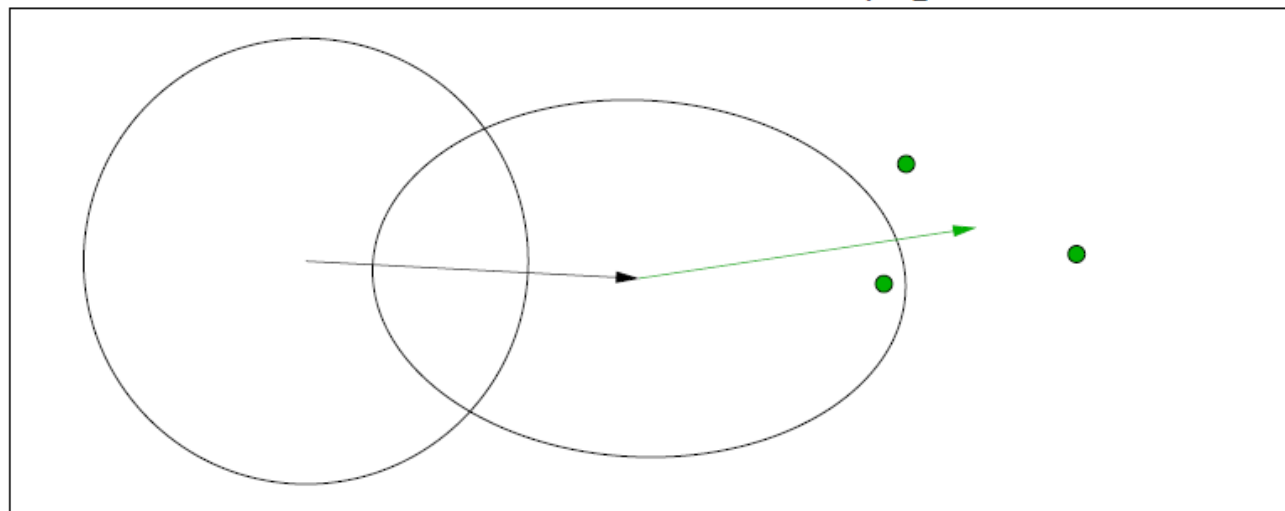
from [Auger, p. 41]

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}

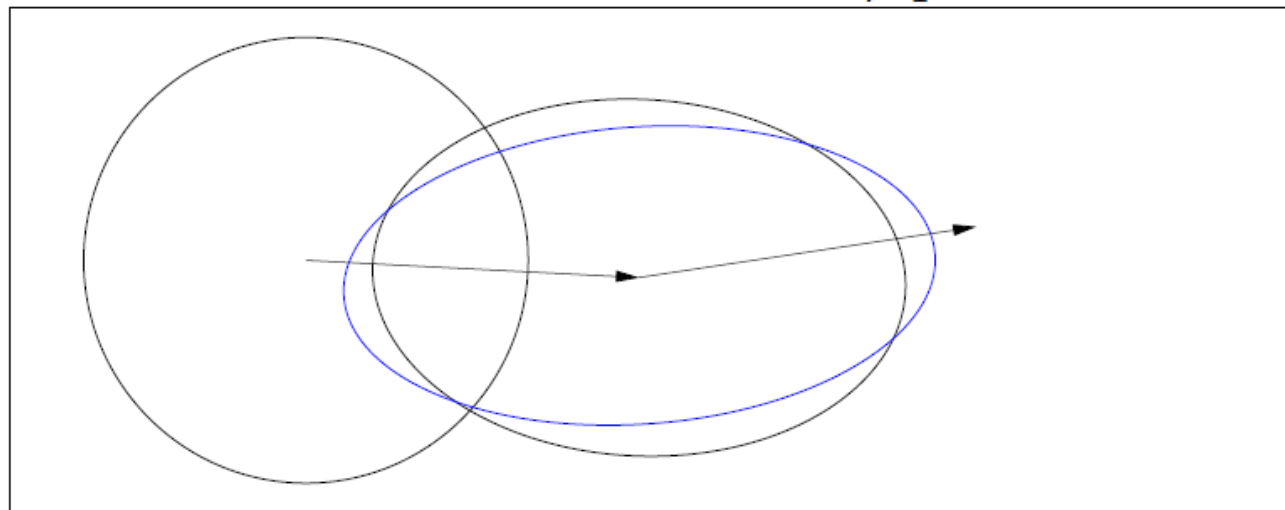
from [Auger, p. 41]

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 ,

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

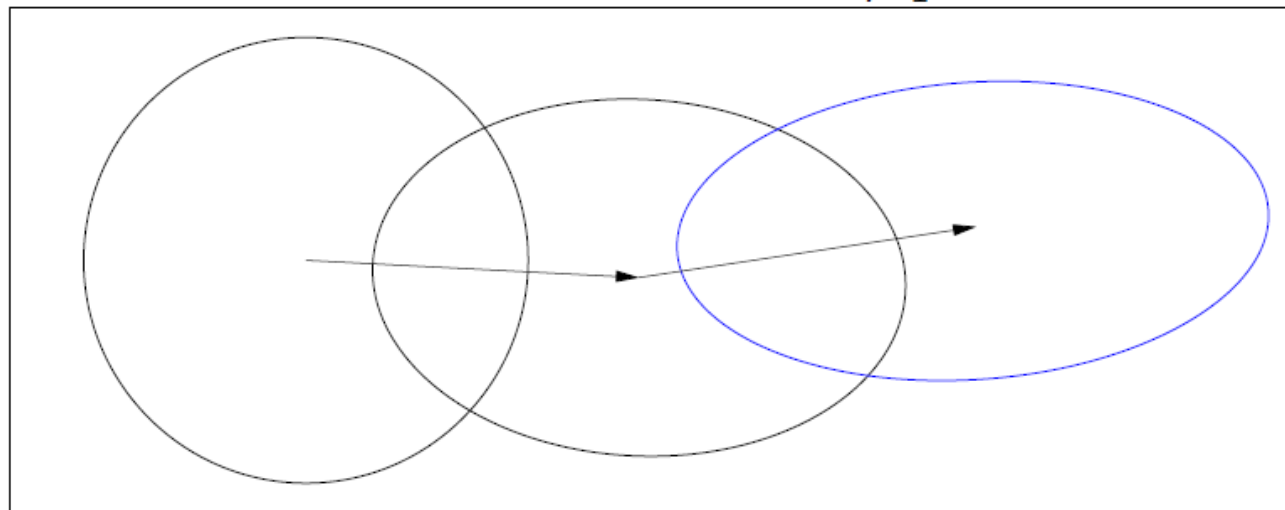
from [Auger, p. 41]

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

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

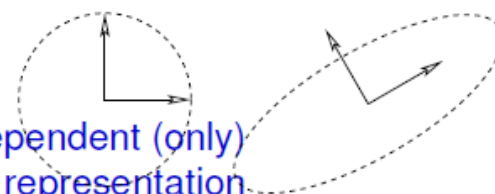
covariance matrix adaptation

- 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

from [Hansen, p. 71]

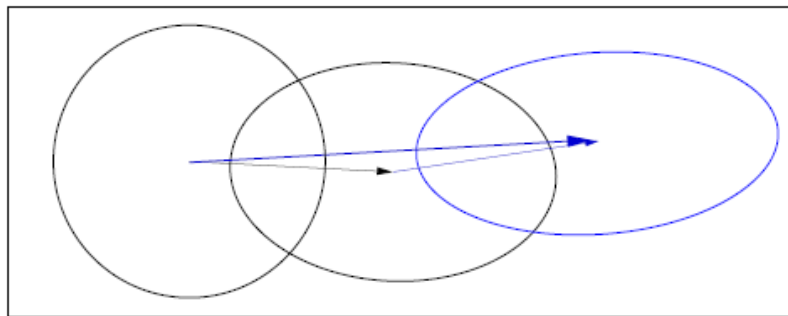
Evolution Path

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.

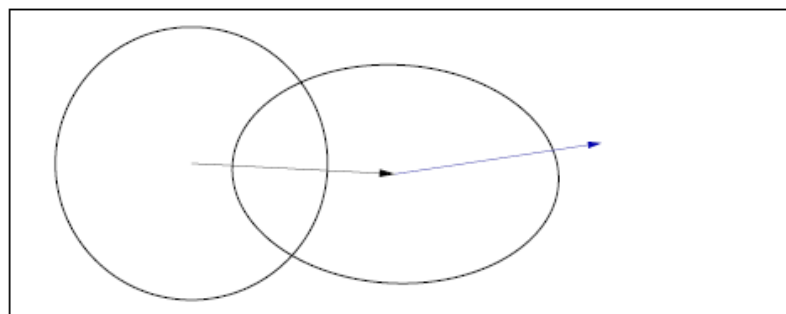
from [Auger, p. 44]

Utilizing the Evolution Path

Cumulation

Utilizing the Evolution Path

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



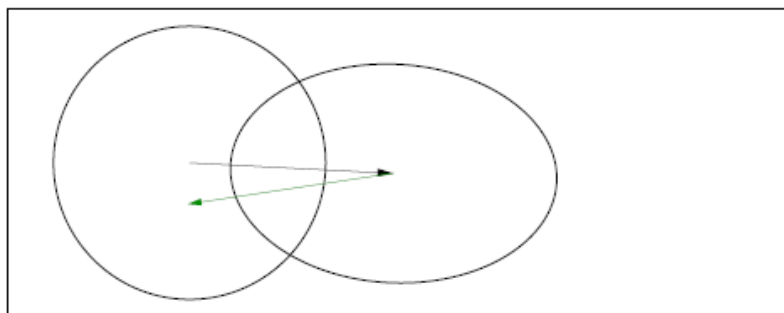
from [Auger, p. 45]

Utilizing the Evolution Path

Cumulation

Utilizing the Evolution Path

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



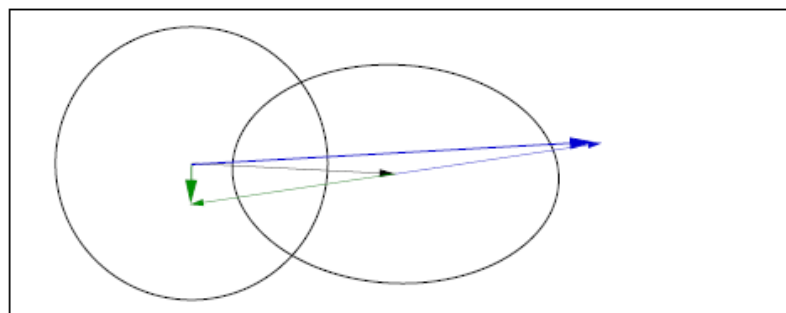
from [Auger, p. 45

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

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

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

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.

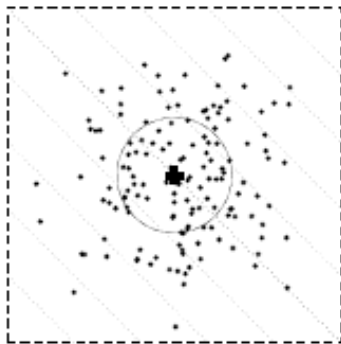
The rank- μ update then reads

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

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

from [Auger, p. 47]

Illustration of Rank- μ Update



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

sampling of

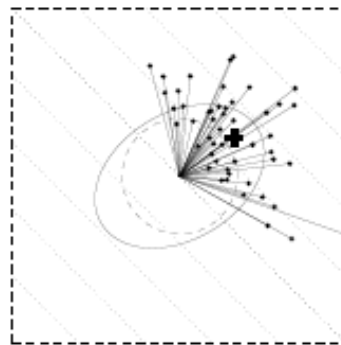
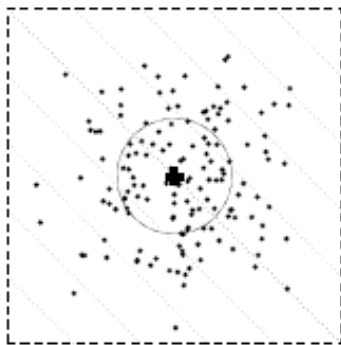
$\lambda = 150$ solutions

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

$$\sigma = 1$$

from [Auger, p. 48]

Illustration of Rank- μ Update



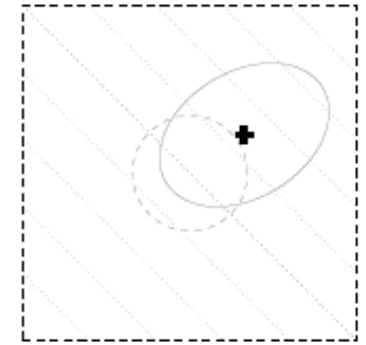
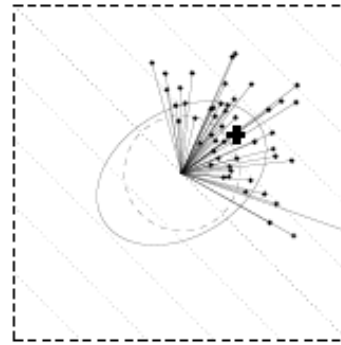
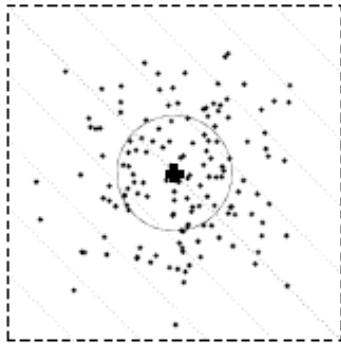
$$x_i = m + \sigma y_i, \quad y_i \sim \mathcal{N}(0, \mathbf{C}) \quad \begin{array}{l} \mathbf{C}_\mu = \frac{1}{\mu} \sum y_{i:\lambda} y_{i:\lambda}^T \\ \mathbf{C} \leftarrow \frac{1}{(\mu - 1) \times \mathbf{C} + 1 \times \mathbf{C}_\mu} \end{array}$$

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
 $c_{\text{cov}} = 1$

from [Auger, p. 48]

Illustration of Rank- μ Update



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

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

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

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

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

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

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: $c_c \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 1$, and $w_{i=1 \dots \lambda}$ such that $\mu_w =$

Promised:

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

$\frac{\mu_w}{n}$,

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 σ

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

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 σ

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

The CMA-ES

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

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

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

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 (0)

What did we want to achieve?

- reduce any convex-quadratic function

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

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

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

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

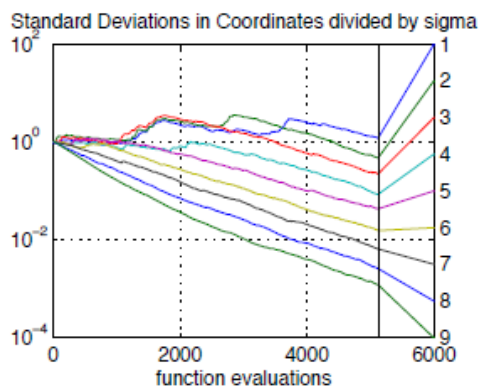
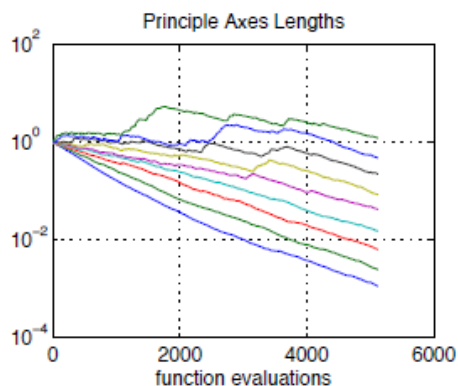
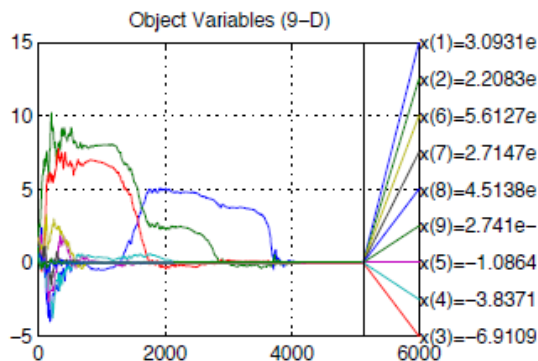
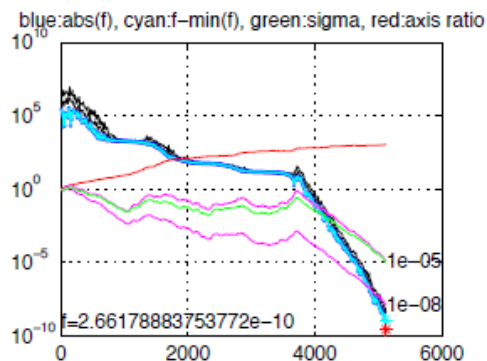
in a stochastic sense

from [Hansen, p. 91]

Experimentum Crucis with CMA-ES

Experimentum Crucis (1)

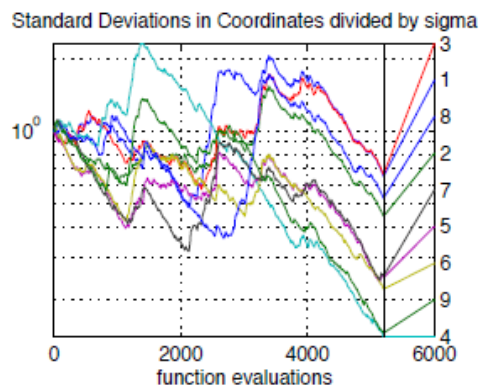
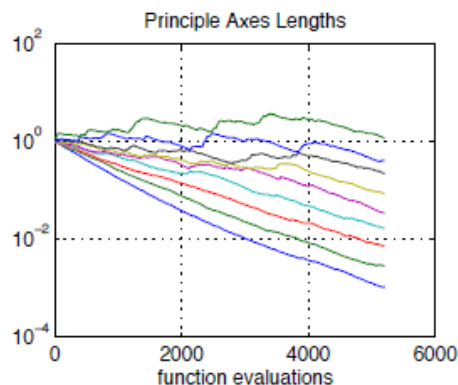
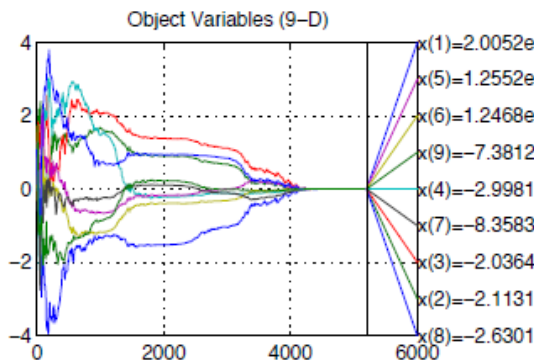
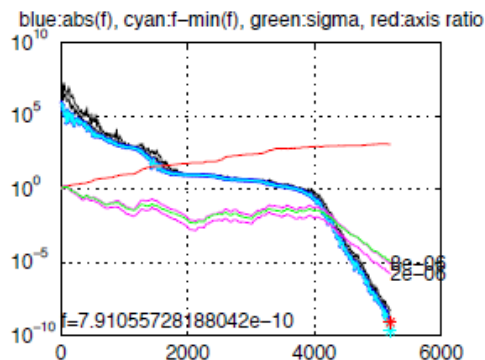
f convex quadratic, separable



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

Experimentum Crucis (2)

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



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

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

from [Hansen, p. 93]

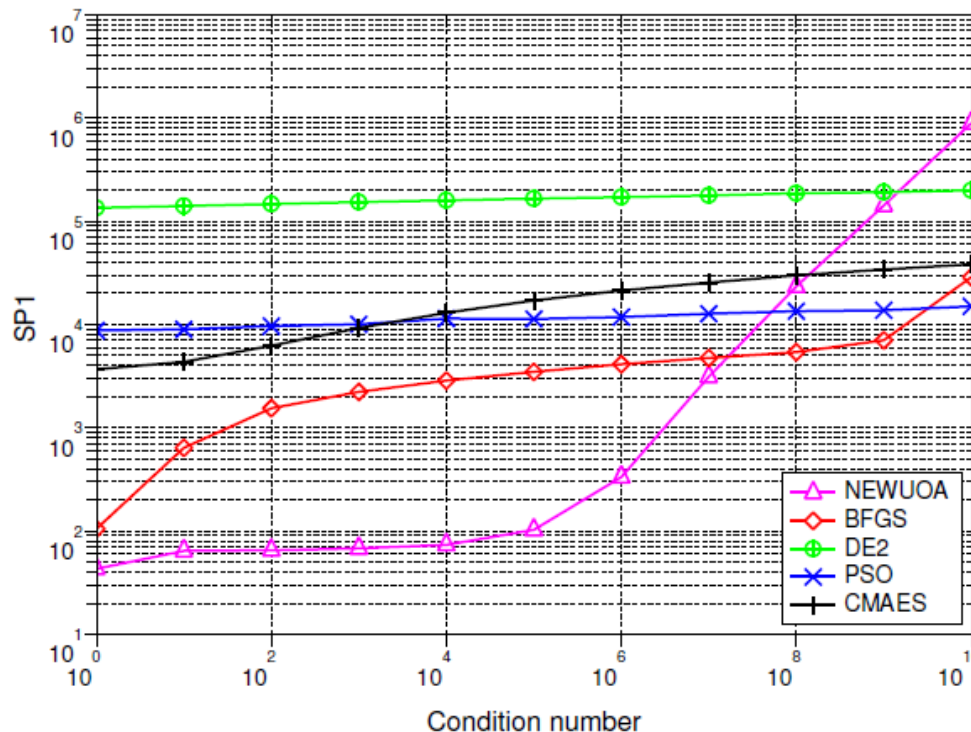
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f convex quadratic, separable with varying condition number 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)

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

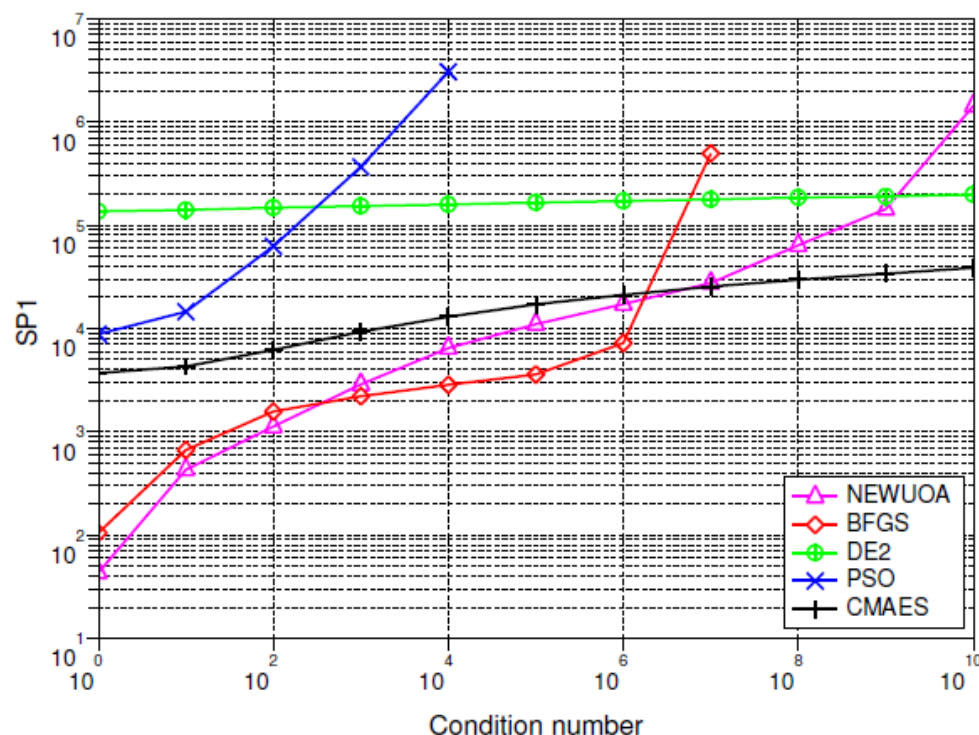
Influence of Condition Number + Invariance

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)

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

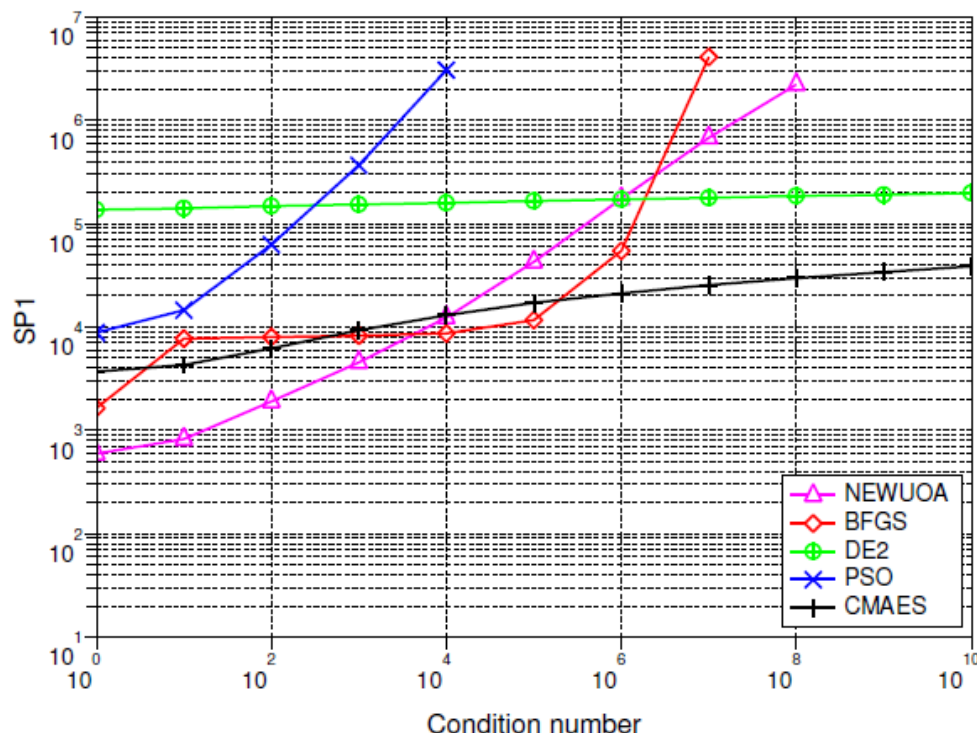
Influence of Condition Number + Invariance

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)

$f(x) = g(x^T H x)$ 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 evaluations¹⁶ to reach the target function value of $g^{-1}(10^{-9})$

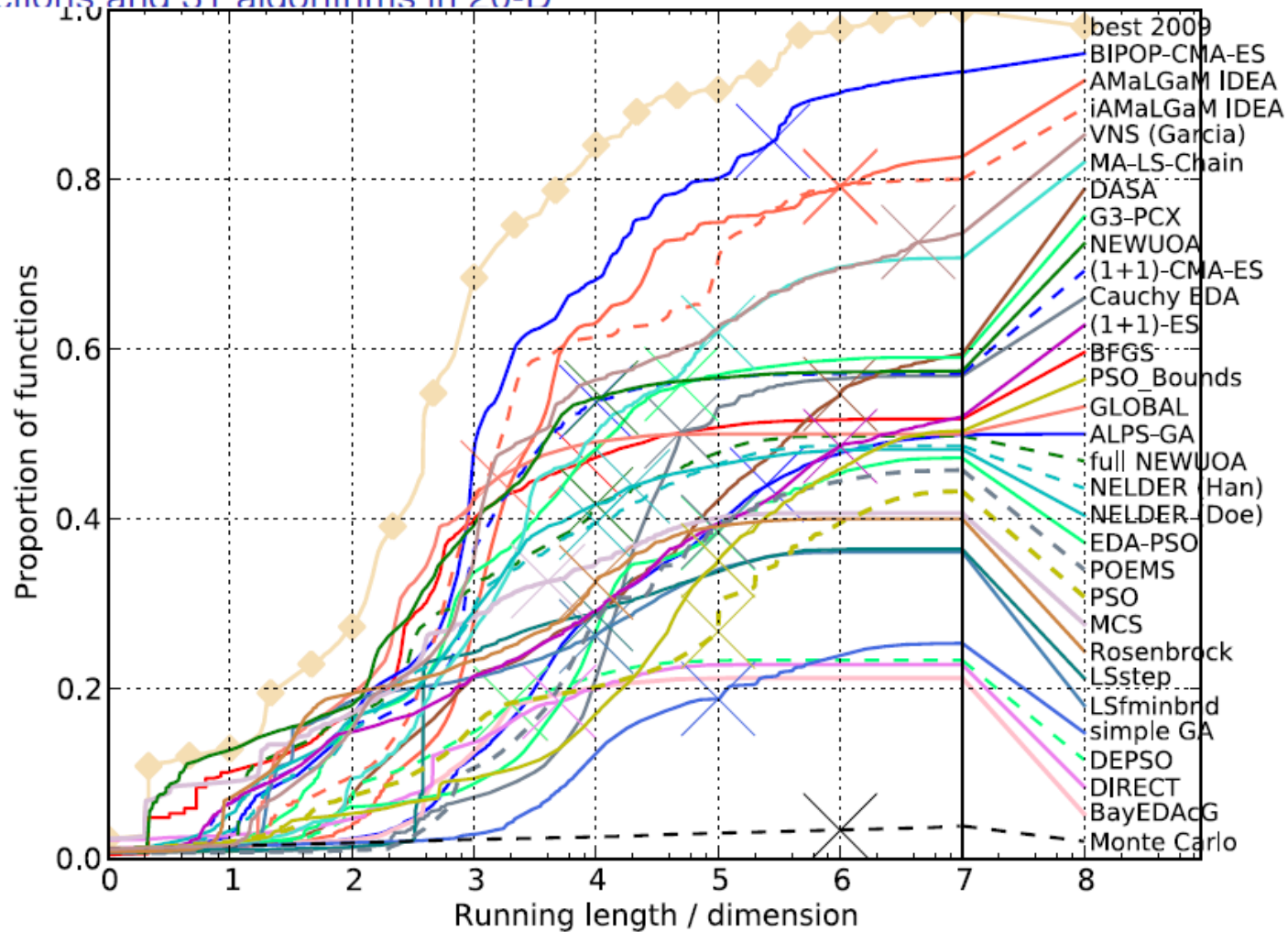
¹⁶ Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

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

Bayesian Optimization

yet another derivative-free approach (but now with surrogates)

Surrogate-Assisted Optimization Idea

When objective function is expensive:

- replace objective function (sometimes) with a (learned) model of it

Questions:

- Which model?
- How to match the model with true evaluations (i.e. how to learn)?
- When to evaluate the objective function, when the model?
- In this context, most algorithms include various heuristic parts...

A Basic Surrogate-Assisted Algorithm Framework

Initial (space-filling) experimental design:

evaluate n_0 points x_1, \dots, x_{n_0} in the search space as $f(x_1), \dots, f(x_{n_0})$ to “get a first idea of the landscape”

While happy do:

- Learn surrogate model \mathcal{M} of the true objective function based on (a subset of) all so-far-evaluated search points
- Find (an approximation of) the optimum of the so-called acquisition function $A_{\mathcal{M}}(x)$, denoted x_n
- Evaluate x_n on the true objective function f

The Framework with Quadratic Model: NEWUOA

(FULL-)NEWUOA by M. J. D. Powell:

- **Initial design** ($t = 0$): $(n + 1)(n + 2)/2$ points around x_0 along coordinate axes
- **Surrogate model**: quadratic, interpolating the current quadratically many points
- **Acquisition function**: minimum $x_{t+1} = x_t + d$ of quadratic model within a trust region (with $\|d\|$ smaller than *adapted* trust region radius Δ)
- Evaluate new point x_{t+1} and replace x_t if better, replace interpolation point, furthest to x_{t+1} , with x_{t+1}
- Many more details to make the algorithm work well in practice
- In particular, often using only $2n + 1$ interpolation points

Powell, M. J. D. (November 2004). The NEWUOA software for unconstrained optimization without derivatives (Report). Department of Applied Mathematics and Theoretical Physics, Cambridge University. DAMTP 2004/NA05.

Powell, M. J. D. (2004). "Least Frobenius norm updating of quadratic models that satisfy interpolation conditions". *Mathematical Programming* 100: 183–215. doi:10.1007/s10107-003-0490-7.

EGO: Gaussian Processes as Probabilistic Model

Bayesian Optimization (Mockus, 1975):

- Model = Gaussian Process
- Gaussian Process aka Kriging: a probabilistic model (see next slide)
- Became popular especially under the name Efficient Global Optimization (EGO, Jones et al., 1998)
- Acquisition function often a trade-off between exploration and exploitation, most common is Expected Improvement

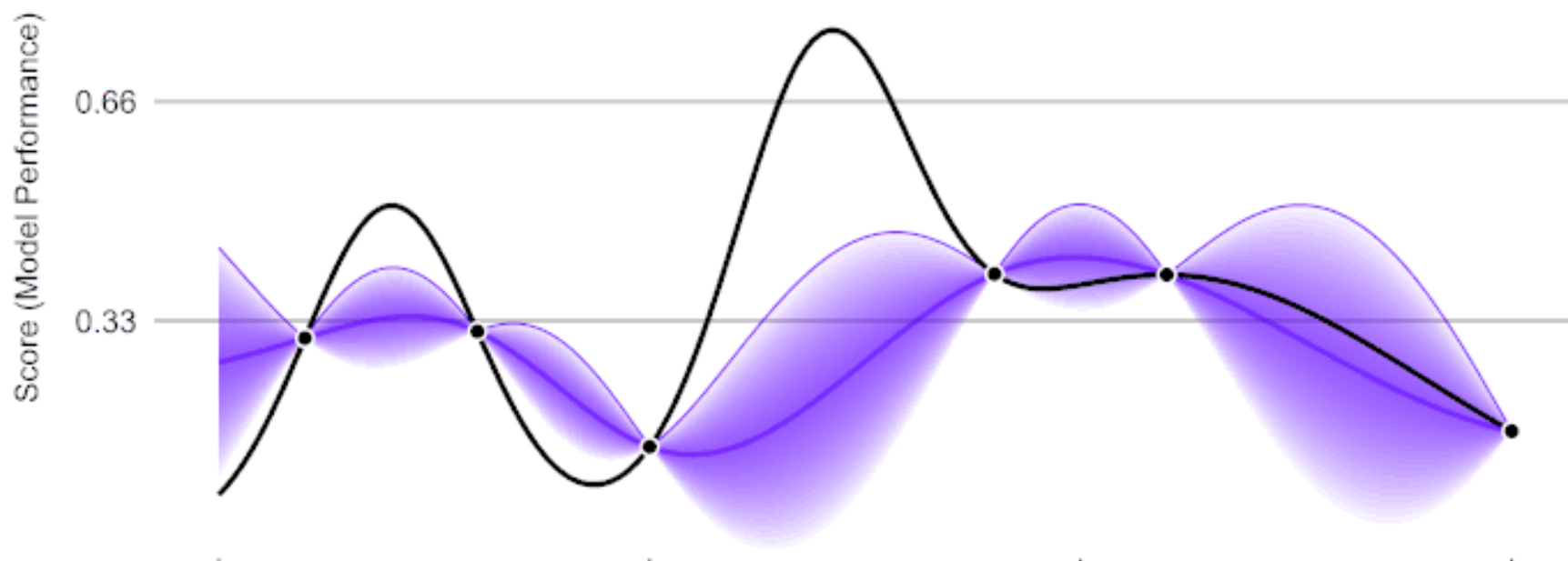
for details, see also <https://arxiv.org/pdf/1807.02811.pdf>

Mockus, J. (1975). On Bayesian methods for seeking the extremum. In Optimization Techniques IFIP Technical Conference, pages 400–404. Springer

Jones, D. R., Schonlau, M., and Welch, W. J. (1998). Efficient global optimization of expensive black-box functions. Journal of Global Optimization, 13(4):455–492.

Gaussian Processes

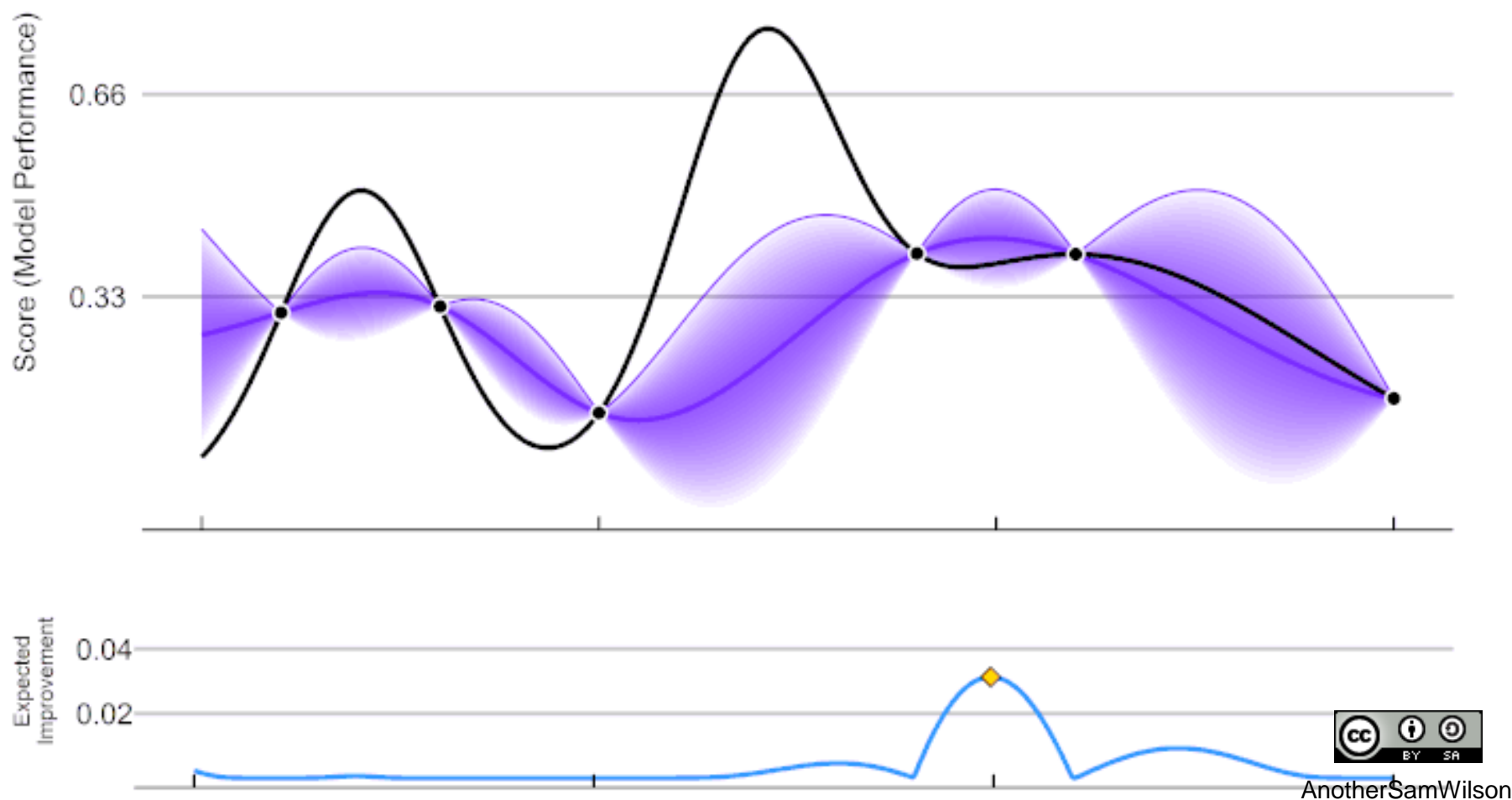
ParBayesianOptimization in Action (Round 1)



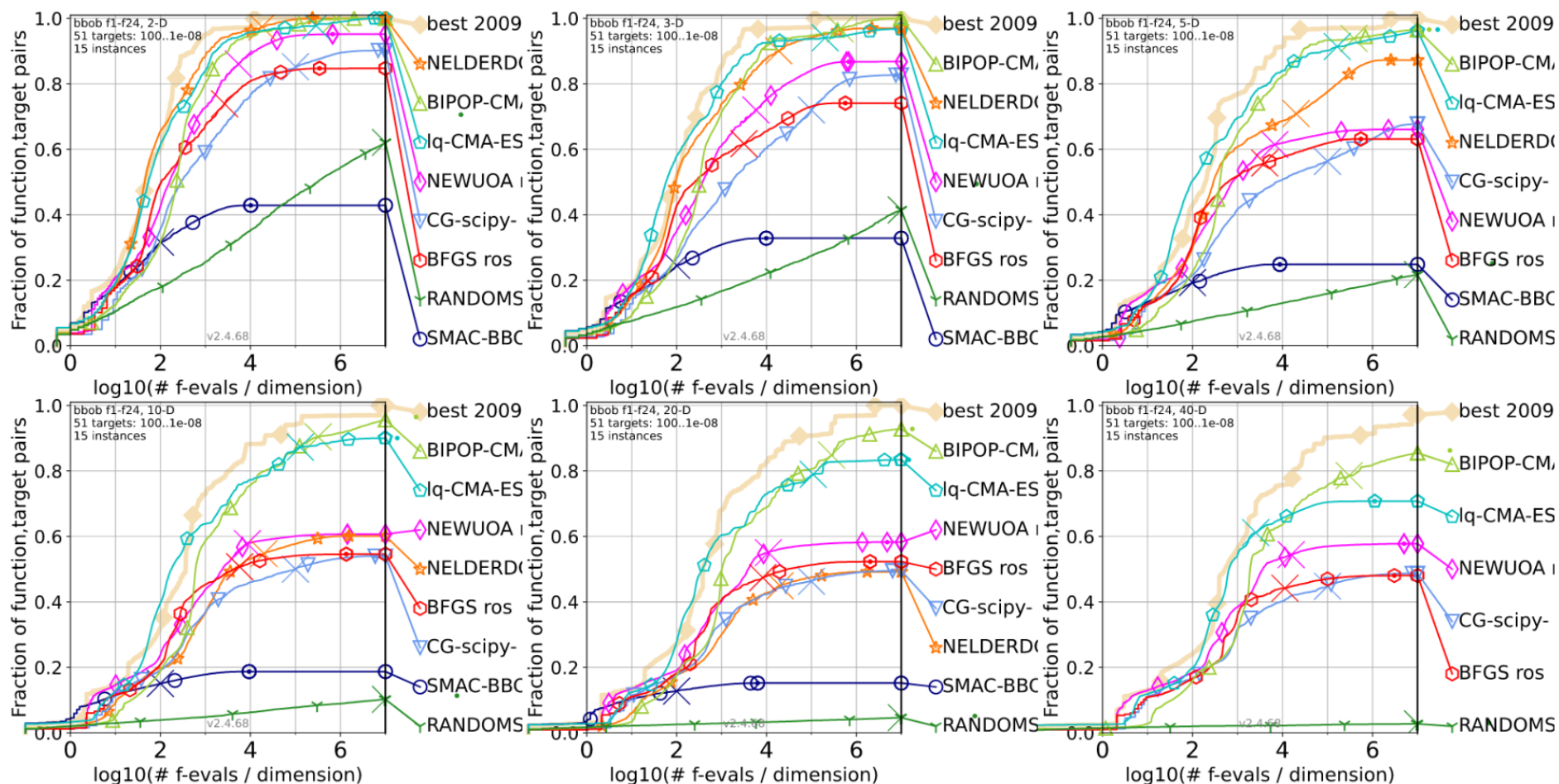
AnotherSamWilson

Gaussian Processes

ParBayesianOptimization in Action (Round 1)



Preview: Comparison of Algorithms



addendum:

Stochastic Gradient Descent

Reminder Gradient Descent

General principle

- ① choose an initial point x_0 , set $t = 0$
- ② while not happy
 - choose the negative gradient as descent direction: $-\nabla f$
 - line search:
 - choose a step size $\sigma_t > 0$
 - set $x_{t+1} = x_t - \sigma_t \nabla f(x_t)$
 - set $t = t + 1$

Disadvantage of Basic Gradient Descent

When optimizing weights of a neural net in deep learning:

- we optimize the loss function $\sum_{i=1}^m (y_i - \hat{y}_i)$
 y_i : from training data, \hat{y}_i : neural network output, depending on weight vector \mathbf{w}
- not only \mathbf{w} is of high dimension, but also m is large
- consequence: computing gradient $\nabla f(\mathbf{w})$ is costly
- Why? → Have to compute many partial derivatives

Idea:

- Compute only an approximation of $\nabla f(\mathbf{w})$ with respect to *some* data
- More generally applicable when $f(x) = \sum_{i=1}^m f_i(x)$

Stochastic Gradient Descent (SGD)

Idea:

Compute only an approximation of $\nabla f(\mathbf{x})$ with respect to *some subfunctions* of $f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x})$

Instead of $\mathbf{x}_{t+1} = \mathbf{x}_t - \sigma_t \nabla f(\mathbf{x}_t)$, we update the current iterate only with respect to one (random) subfunction f_i :

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \sigma_t \nabla f_i(\mathbf{x}_t)$$

Batches, Mini-Batches and Epochs

- To make sure, all subfunctions are used equally, we typically choose a (random) permutation π and use each subfunction once in this order before to redo this step
- Such step is called an *epoch*
- Using single subfunctions is fast, but also highly stochastic
consequence: practical convergence is slower
- Hence, we shall use batches of subfunctions at a time:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \sigma_t \nabla f_{\mathcal{B}_j}(\mathbf{x}_t)$$

with $f_{\mathcal{B}_j}(\mathbf{x}_t) = \sum_{i \in \mathcal{B}_j} f_i(\mathbf{x})$ and $\mathcal{B}_j \subset \{1, \dots, m\}$ the corresponding indices in the j th batch

Stochastic Gradient Descent

Typically applied with a constant step size in the context of Deep Learning

Adam: Adaptive Moment Adaptation

- decaying average of past gradients $g_t(x) = \nabla f_{\mathcal{B}_j}(x_t)$, also called “momentum”:

$$m_t(x) = \beta_1 m_{t-1} + (1 - \beta_1) g_t(x)$$

- decaying average of past squared gradients

$$v_t(x) = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2(x)$$

- to decrease bias towards 0, we use $\widehat{m}_t = \frac{m_t}{1 - \beta_1^t}$ and $\widehat{v}_t = \frac{v_t}{1 - \beta_2^t}$
- and then $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\sigma_t}{\sqrt{\widehat{v}_t(x)} + \epsilon} \widehat{m}_t(x)$