

# Introduction to Randomized Continuous Optimization

**Youhei Akimoto<sup>1</sup> & Anne Auger<sup>2</sup> & Nikolaus Hansen<sup>2</sup>**  
1. Shinshu University, Nagano, Japan  
2. Inria, Research Centre Saclay, France

y\_akimoto@shinshu-u.ac.jp  
anne.auger@inria.fr  
nikolaus.hansen@inria.fr

<http://www.sigevo.org/gecco-2016/>



Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s). Copyright is held by the author(s).  
GECCO '16 Companion, July 20-24, 2016, Denver, CO, USA.  
ACM 978-1-4503-4323-7/16/07.  
<http://dx.doi.org/10.1145/2936981.2926980>

1

We are happy to answer questions at any time.

2

## Overview

### 1 Problem Statement

Continuous Black-Box Optimization  
Typical Difficulties

### 2 Stochastic Black-Box Algorithms

General Template  
Invariance  
Comparisons of a few DFOs

### 3 Zoom on Evolution Strategies

Step-size Adaptation  
Covariance Matrix Adaptation

### 4 Evaluating Black-Box Algorithms

Displaying results and visualization  
Statistics  
Average Runtime  
Empirical Cumulative Distribution Function (ECDF)

3

## Problem Statement

### Continuous Domain Search/Optimization

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

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

- Black Box** scenario (direct search scenario)



- gradients are not available or not useful
- problem domain specific knowledge is used only within the black box, e.g. within an appropriate encoding
- Search **costs**: number of function evaluations

4

## Problem Statement

Continuous Domain Search/Optimization

- Goal
  - ▶ fast convergence to the global optimum
  - ▶ solution  $x$  with **small function value**  $f(x)$  with **least search cost**  
 ... or to a robust solution  $x$   
 there are two conflicting objectives
  
- Typical Examples
  - ▶ shape optimization (e.g. using CFD) curve fitting, airfoils
  - ▶ model calibration biological, physical
  - ▶ parameter calibration controller, plants, images
  
- Problems
  - ▶ exhaustive search is infeasible
  - ▶ naive random search takes too long
  - ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

## Problem Statement

Continuous Domain Search/Optimization

- Goal
  - ▶ fast convergence to the global optimum
  - ▶ solution  $x$  with **small function value**  $f(x)$  with **least search cost**  
 ... or to a robust solution  $x$   
 there are two conflicting objectives
  
- Typical Examples
  - ▶ shape optimization (e.g. using CFD) curve fitting, airfoils
  - ▶ model calibration biological, physical
  - ▶ parameter calibration controller, plants, images
  
- Problems
  - ▶ exhaustive search is infeasible
  - ▶ naive random search takes too long
  - ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

## Problem Statement

Continuous Domain Search/Optimization

- Goal
  - ▶ fast convergence to the global optimum
  - ▶ solution  $x$  with **small function value**  $f(x)$  with **least search cost**  
 ... or to a robust solution  $x$   
 there are two conflicting objectives
  
- Typical Examples
  - ▶ shape optimization (e.g. using CFD) curve fitting, airfoils
  - ▶ model calibration biological, physical
  - ▶ parameter calibration controller, plants, images
  
- Problems
  - ▶ exhaustive search is infeasible
  - ▶ naive random search takes too long
  - ▶ deterministic search is not successful / takes too long

Approach: stochastic search, Evolutionary Algorithms

## Objective Function Properties

We assume  $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  to be *non-linear*, *non-separable* and to have at least moderate dimensionality, say  $n \ll 10$ .

Additionally,  $f$  can be

- non-convex there are possibly many local optima
- multimodal
- non-smooth derivatives do not exist
- discontinuous, plateaus
- ill-conditioned
- noisy
- ...

Goal : cope with any of these function properties  
 they are related to real-world problems

## Objective Function Properties

We assume  $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$  to be *non-linear, non-separable* and to have at least moderate dimensionality, say  $n \ll 10$ .

Additionally,  $f$  can be

- non-convex
- multimodal
- non-smooth
- discontinuous, plateaus
- ill-conditioned
- noisy
- ...

there are possibly many local optima

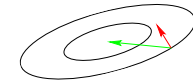
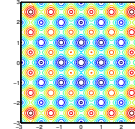
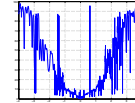
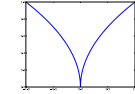
derivatives do not exist

**Goal:** cope with any of these function properties  
they are related to real-world problems

## What Makes a Function Difficult to Solve?

Why stochastic search?

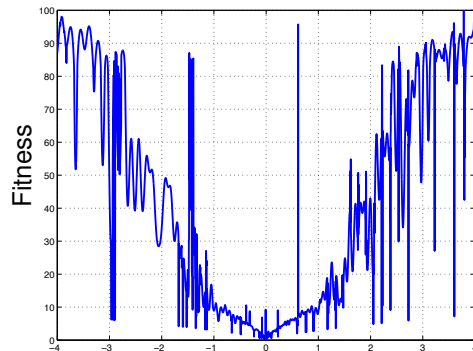
- non-linear, non-quadratic, non-convex  
on linear and quadratic functions much better search policies are available
- ruggedness  
non-smooth, discontinuous, multimodal, and/or noisy function
- dimensionality (size of search space)  
(considerably) larger than three
- non-separability  
dependencies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



gradient direction Newton direction

## Ruggedness

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



cut from a 5-D example, (easily) solvable with evolution strategies

## Curse of Dimensionality

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

Example: Consider placing 20 points equally spaced onto the interval  $[0, 1]$ . Now consider the 10-dimensional space  $[0, 1]^{10}$ . To get **similar coverage** in terms of distance between adjacent points requires  $20^{10} \approx 10^{13}$  points. 20 points appear now as isolated points in a vast empty space.

Remark: **distance measures** break down in higher dimensionalities (the central limit theorem kicks in)

Consequence: a **search policy** that is valuable in small dimensions **might be useless** in moderate or large dimensional search spaces. Example: exhaustive search.

## Separable Problems

### Definition (Separable Problem)

A function  $f$  is separable if

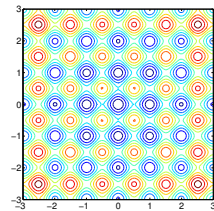
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left( \arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

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

### Example: Additively decomposable functions

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

Rastrigin function



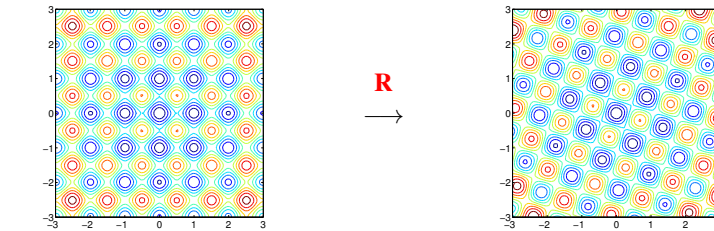
## Non-Separable Problems

Building a non-separable problem from a separable one <sup>(1,2)</sup>

### Rotating the coordinate system

- $f : \mathbf{x} \mapsto f(\mathbf{x})$  separable
- $f : \mathbf{x} \mapsto f(\mathbf{R}\mathbf{x})$  non-separable

$\mathbf{R}$  rotation matrix



<sup>1</sup>Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann  
<sup>2</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

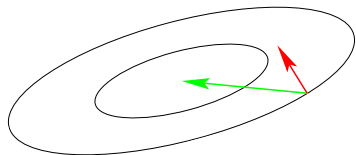
## Ill-Conditioned Problems

Curvature of level sets

Consider the convex-quadratic function

$$f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x} - \mathbf{x}^*) = \frac{1}{2} \sum_i h_{i,i} (x_i - x_i^*)^2 + \frac{1}{2} \sum_{i \neq j} h_{i,j} (x_i - x_i^*)(x_j - x_j^*)$$

$\mathbf{H}$  is Hessian matrix of  $f$  and symmetric positive definite



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

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

Ill-conditioning means **squeezed level sets** (high curvature).  
 Condition number equals nine here. Condition numbers up to  $10^{10}$   
 are not unusual in real world problems.

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

## What Makes a Function Difficult to Solve?

... and what can be done

The Problem	Possible Approaches
Dimensionality	exploiting the problem structure separability, locality/neighborhood, encoding
Ill-conditioning	second order approach changes the neighborhood metric
Ruggedness	<b>non-local</b> policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed
	<b>population-based</b> method, stochastic, non-elitistic recombination operator serves as repair mechanism
	restarts

... metaphors

## Metaphors

Evolutionary Computation Optimization/Nonlinear Programming

individual, offspring, parent	↔	candidate solution
		decision variables
		design variables
		object variables
population	↔	set of candidate solutions
fitness function	↔	objective function
		loss function
		cost function
		error function
generation	↔	iteration

... methods: ESs

## Landscape of Continuous Black-Box Optimization

### Deterministic algorithms

- Quasi-Newton with estimation of gradient (BFGS) [Broyden et al. 1970]
- Simplex downhill [Nelder & Mead 1965]
- Pattern search [Hooke and Jeeves 1961]
- Trust-region methods (NEWUOA, BOBYQA) [Powell 2006, 2009]

### Stochastic (randomized) search methods

- Evolutionary Algorithms (continuous domain)
  - Differential Evolution [Storn & Price 1997]
  - Particle Swarm Optimization [Kennedy & Eberhart 1995]
  - Evolution Strategies, CMA-ES** [Rechenberg 1965, Hansen & Ostermeier 2001]
  - Estimation of Distribution Algorithms (EDAs) [Larrañaga, Lozano, 2002]
  - Cross Entropy Method (same as EDA) [Rubinstein, Kroese, 2004]
  - Genetic Algorithms [Holland 1975, Goldberg 1989]
- Simulated annealing [Kirkpatrick et al. 1983]
- Simultaneous perturbation stochastic approximation (SPSA) [Spall 2000]

## Overview

- 1 Problem Statement
  - Continuous Black-Box Optimization
  - Typical Difficulties
- 2 Stochastic Black-Box Algorithms
  - General Template
  - Invariance
  - Comparisons of a few DFOs
- 3 Zoom on Evolution Strategies
  - Step-size Adaptation
  - Covariance Matrix Adaptation
- 4 Evaluating Black-Box Algorithms
  - Displaying results and visualization
  - Statistics
  - Average Runtime
  - Empirical Cumulative Distribution Function (ECDF)

## Stochastic Search

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

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

While not terminate

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

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

deterministic algorithms are covered as well

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

Natural template for (incremental) *Estimation of Distribution Algorithms*

## Examples

- 1 Estimation of Distribution Algorithms
- 2 Evolution Strategies
- 3 Differential Evolution
- 4 Particle Swarm Optimization

*all those methods are comparison based*

## Evolution Strategies

New search points are sampled normally distributed

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

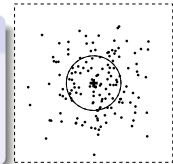
as perturbations of  $m$ , where  $x_i, m \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ ,  $C \in \mathbb{R}^{n \times n}$

where

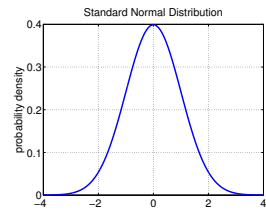
- the **mean** vector  $m \in \mathbb{R}^n$  represents the favorite solution
- the so-called **step-size**  $\sigma \in \mathbb{R}_+$  controls the *step length*
- the **covariance matrix**  $C \in \mathbb{R}^{n \times n}$  determines the **shape** of the distribution ellipsoid

here, all new points are sampled with the same parameters

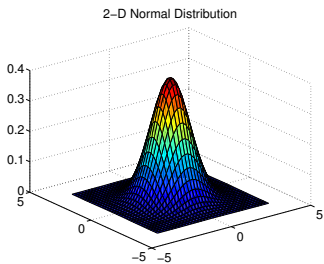
The question remains how to update  $m$ ,  $C$ , and  $\sigma$ .



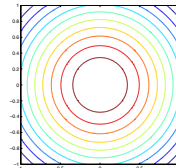
## Normal Distribution



probability density of the 1-D standard normal distribution



probability density of a 2-D normal distribution

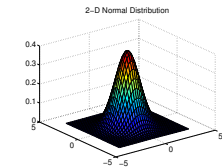


## The Multi-Variate ( $n$ -Dimensional) Normal Distribution

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

The **mean** value  $m$

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

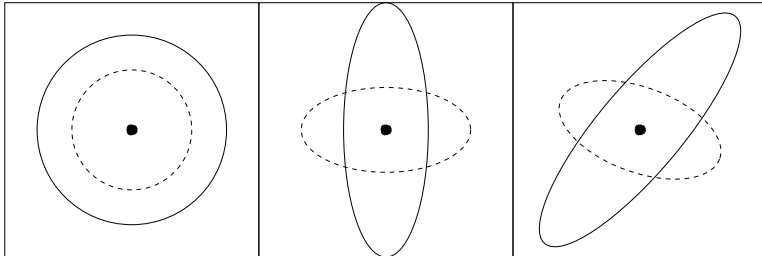


The **covariance matrix**  $C$

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

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

Lines of Equal Density



$\mathcal{N}(m, \sigma^2 \mathbf{I}) \sim m + \sigma \mathcal{N}(\mathbf{0}, \mathbf{I})$      $\mathcal{N}(m, \mathbf{D}^2) \sim m + \mathbf{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$      $\mathcal{N}(m, \mathbf{C}) \sim m + \mathbf{C}^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$   
**one degree of freedom**  $\sigma$  components are independent standard normally distributed  
 **$n$  degrees of freedom** components are independent, scaled  
 **$(n^2 + n)/2$  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  $\mathbf{A}$ .

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

Non-elitist selection and intermediate (weighted) recombination

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

Let  $x_{i:\lambda}$  the  $i$ -th ranked solution point, such that  $f(x_{1:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda})$ .  
 The new mean reads

$$m \leftarrow \sum_{i=1}^{\mu} w_i x_{i:\lambda} = m + \sigma \underbrace{\sum_{i=1}^{\mu} w_i y_{i:\lambda}}_{=: y_w}$$

where

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

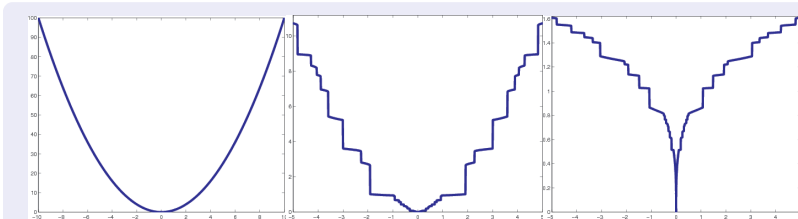
The best  $\mu$  points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

### Invariance Under Monotonically Increasing Functions

#### Rank-based algorithms

Update of all parameters uses only the ranks

$$f(x_{1:\lambda}) \leq f(x_{2:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda})$$



$$g(f(x_{1:\lambda})) \leq g(f(x_{2:\lambda})) \leq \dots \leq g(f(x_{\lambda:\lambda})) \quad \forall g$$

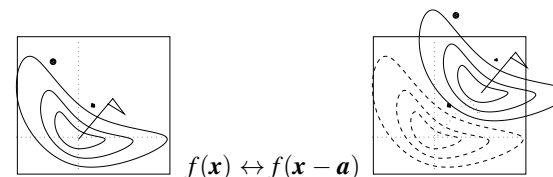
$g$  is strictly monotonically increasing  
 $g$  preserves ranks

<sup>3</sup>Whitley 1989. The GENITOR algorithm and selection pressure: Why rank-based allocation of reproductive trials is best, ICGA

### Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



$$f(x) \leftrightarrow f(x - a)$$

#### Identical behavior on $f$ and $f_a$

$$f: x \mapsto f(x), \quad x^{(t=0)} = x_0$$

$$f_a: x \mapsto f(x - a), \quad x^{(t=0)} = x_0 + a$$

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

Evolution Strategies (ES) Invariance

### Invariance Under Rigid Search Space Transformations

$f = h_{\text{Rast}}$   $f$ -level sets in dimension 2  $f = h$

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

29 Anne Auger & Nikolaus Hansen CMA-ES July, 2014 27 / 81

Evolution Strategies (ES) Invariance

### 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)

30 Anne Auger & Nikolaus Hansen CMA-ES July, 2014 27 / 81

Evolution Strategies (ES) Invariance

### 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

31

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

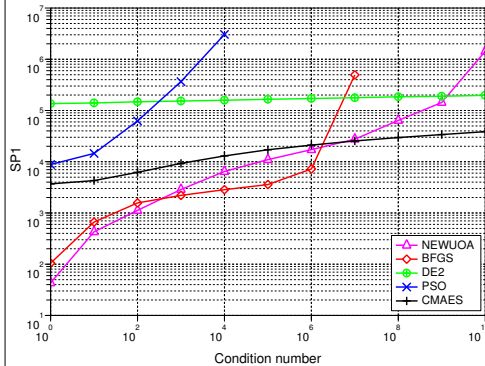
32



## 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)

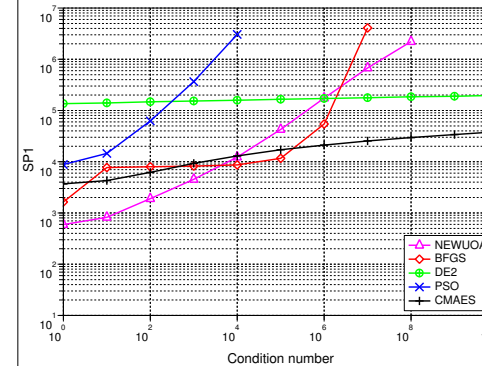
SP1 = average number of objective function evaluations<sup>15</sup> to reach the target function value of  $g^{-1}(10^{-9})$

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

## 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> Auer et al. (2009): Experimental comparisons of derivative free optimization algorithms. SEA

## Overview

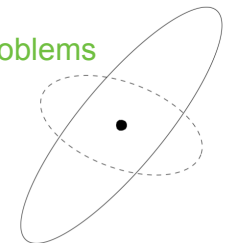
- 1 Problem Statement
  - Continuous Black-Box Optimization
  - Typical Difficulties
- 2 Stochastic Black-Box Algorithms
  - General Template
  - Invariance
  - Comparisons of a few DFOs
- 3 Zoom on Evolution Strategies
  - Step-size Adaptation
  - Covariance Matrix Adaptation
- 4 Evaluating Black-Box Algorithms
  - Displaying results and visualization
  - Statistics
  - Average Runtime
  - Empirical Cumulative Distribution Function (ECDF)

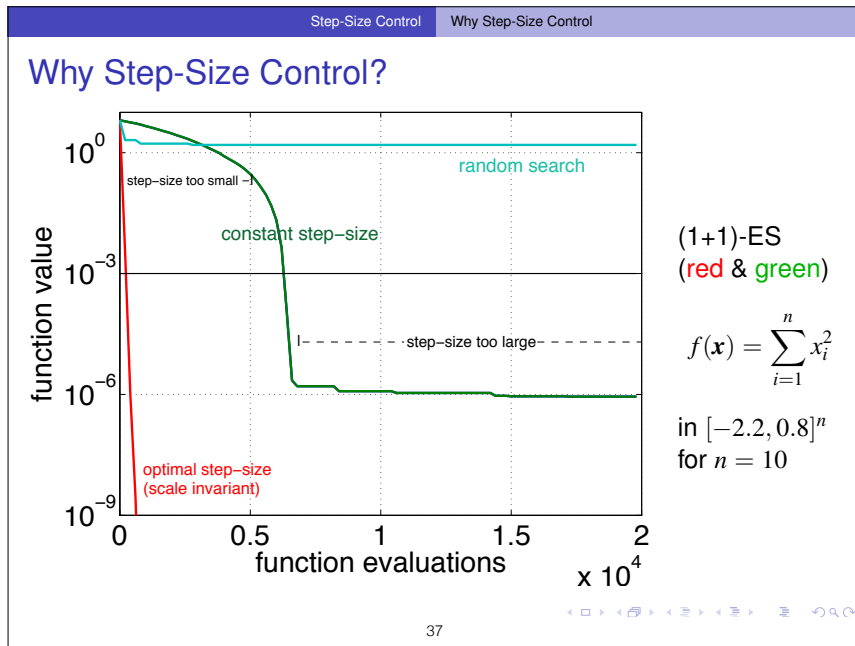
## Zoom on ESs: Objectives

Illustrate **why and how** sampling distribution is controlled

step-size control (overall standard deviation)  
 allows to achieve linear convergence

covariance matrix control  
 allows to solve ill-conditioned problems





Step-Size Control Why Step-Size Control

## Methods for Step-Size Control

- **1/5-th success rule<sup>ab</sup>**, often applied with “+”-selection
  - increase step-size if more than 20% of the new solutions are successful, decrease otherwise
- **$\sigma$ -self-adaptation<sup>c</sup>**, applied with “-”-selection
  - mutation is applied to the step-size and the better, according to the objective function value, is selected
  - simplified “global” self-adaptation
- **path length control<sup>d</sup>** (Cumulative Step-size Adaptation, CSA)<sup>e</sup>
  - self-adaptation derandomized and non-localized

<sup>a</sup>Rechenberg 1973, *Evolutionsstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog  
<sup>b</sup>Schumer and Steiglitz 1968, Adaptive step size random search. *IEEE TAC*  
<sup>c</sup>Schwefel 1981, *Numerical Optimization of Computer Models*, Wiley  
<sup>d</sup>Hansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.* 9(2)  
<sup>e</sup>Ostermeier et al 1994, Step-size adaptation based on non-local use of selection information, *PPSN IV*

38

Step-Size Control Path Length Control (CSA)

## Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

$$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w$$

Measure the length of the *evolution path*  
the pathway of the mean vector  $\mathbf{m}$  in the generation sequence

loosely speaking steps are

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

39

Step-Size Control Path Length Control (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  $c_\sigma \approx 4/n$ ,  $d_\sigma \approx 1$ .

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

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

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

$> 1 \iff \|\mathbf{p}_\sigma\|$  is greater than its expectation

40

## Path Length Control (CSA)

The Equations

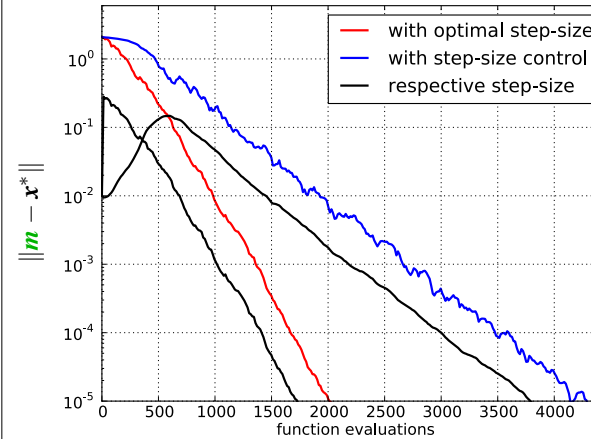
Initialize  $\mathbf{m} \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ , evolution path  $\mathbf{p}_\sigma = \mathbf{0}$ ,  
 set  $c_\sigma \approx 4/n$ ,  $d_\sigma \approx 1$ .

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

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

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

## (5/5, 10)-CSA-ES, default parameters



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

in  $[-0.2, 0.8]^n$   
 for  $n = 30$

## Overview

- 1 Problem Statement
  - Continuous Black-Box Optimization
  - Typical Difficulties
- 2 Stochastic Black-Box Algorithms
  - General Template
  - Invariance
  - Comparisons of a few DFOs
- 3 Zoom on Evolution Strategies
  - Step-size Adaptation
  - Covariance Matrix Adaptation
- 4 Evaluating Black-Box Algorithms
  - Displaying results and visualization
  - Statistics
  - Average Runtime
  - Empirical Cumulative Distribution Function (ECDF)

## Evolution Strategies

Recalling

New search points are sampled normally distributed

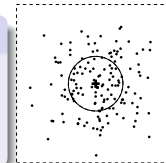
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

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

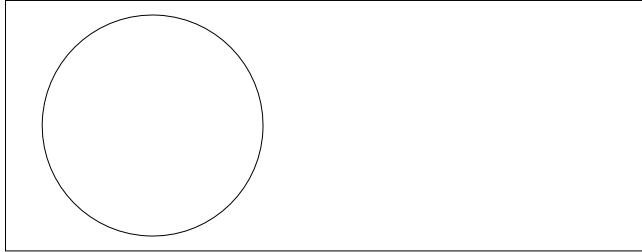
The remaining question is how to update  $\mathbf{C}$ .



### Covariance Matrix Adaptation

Rank-One Update

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



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

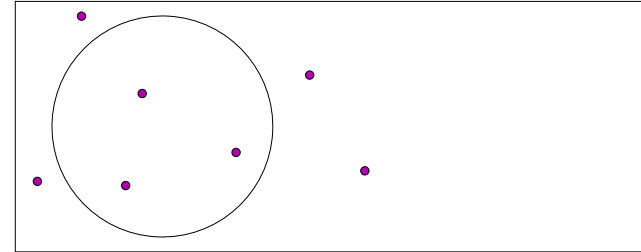
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



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

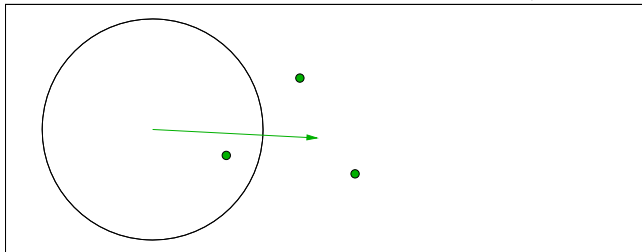
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



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

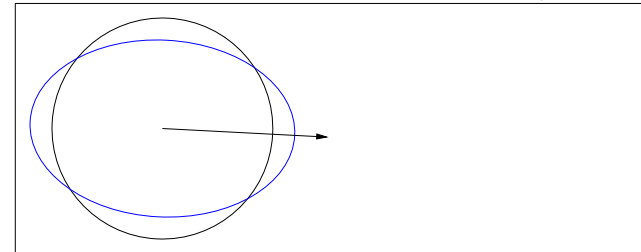
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



mixture of distribution  $\mathbf{C}$  and step  $\mathbf{y}_w$ ,  
 $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$

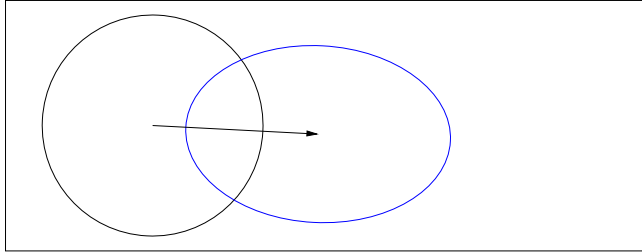
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



new distribution (disregarding  $\sigma$ )

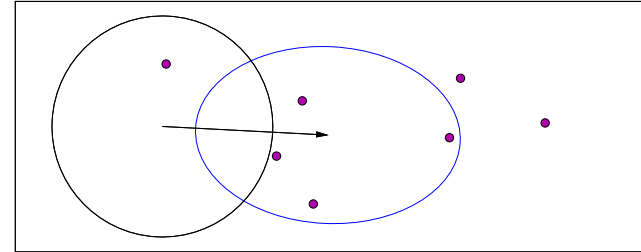
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



new distribution (disregarding  $\sigma$ )

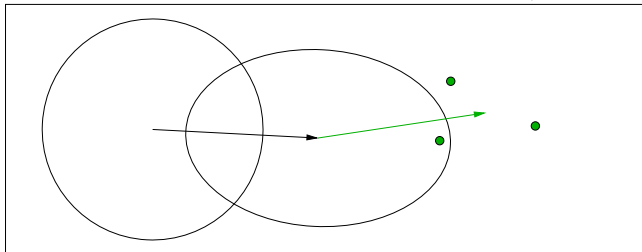
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



movement of the population mean  $\mathbf{m}$

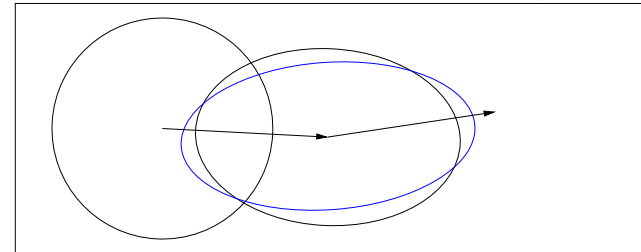
... equations



### Covariance Matrix Adaptation

Rank-One Update

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



mixture of distribution  $\mathbf{C}$  and step  $\mathbf{y}_w$ ,  
 $\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$

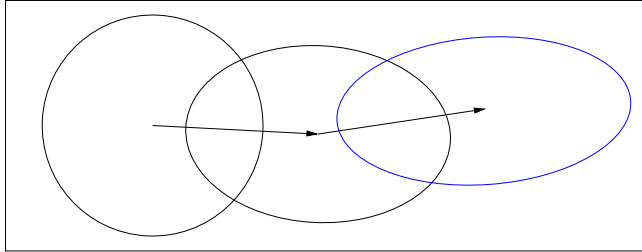
... equations



## 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

another viewpoint: the adaptation **follows a natural gradient**

approximation of the expected fitness

... equations

## Covariance Matrix Adaptation

Rank-One Update

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

While not terminate

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

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

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

The rank-one update has been found independently in several domains<sup>6 7 8 9</sup>

<sup>6</sup> Kjellström&Taxén 1981. Stochastic Optimization in System Design, IEEE TCS  
<sup>7</sup> Hansen&Ostermeier 1996. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation, ICEC  
<sup>8</sup> Ljung 1999. System Identification: Theory for the User  
<sup>9</sup> Haario et al 2001. An adaptive Metropolis algorithm, JSTOR

## 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

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

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

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

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

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

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

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

## Experimentum Crucis (0)

What did we want to achieve?

- reduce any convex-quadratic function

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

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

to the sphere model

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

without use of derivatives

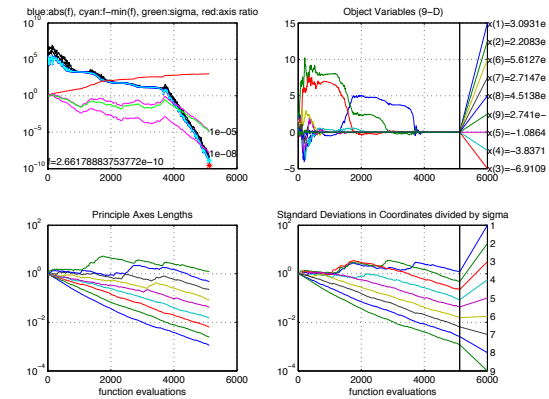
- lines of equal density align with lines of equal fitness

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

in a stochastic sense

## Experimentum Crucis (1)

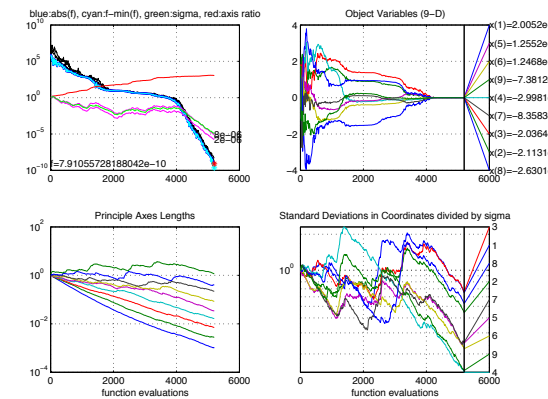
$f$  convex quadratic, separable



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

## Experimentum Crucis (2)

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



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

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

## Overview

- 1 Problem Statement
  - Continuous Black-Box Optimization
  - Typical Difficulties
- 2 Stochastic Black-Box Algorithms
  - General Template
  - Invariance
  - Comparisons of a few DFOs
- 3 Zoom on Evolution Strategies
  - Step-size Adaptation
  - Covariance Matrix Adaptation
- 4 Evaluating Black-Box Algorithms
  - Displaying results and visualization
  - Statistics
  - Average Runtime
  - Empirical Cumulative Distribution Function (ECDF)

## Evaluation of Anytime Black-Box Optimizers

Particularly Randomized Search Algorithms

Randomized optimization is mostly an empirical science

Hence it is crucial to properly conduct numerical experiments and assess performance

to not fool oneself on what our favorite algorithm is good/not good at

in order to not fool others ...

*“The first principle is that you must not fool yourself and you are the easiest person to fool.”*

Richard P. Feynman

## Evaluation of Anytime Black-Box Optimizers

Particularly Randomized Search Algorithms

Evaluation of performance in a **broad sense**

not only be able to say “Algorithm A is better than B”

but

**understand** where and why algorithm work

**quantify** performance

61

## Measuring Performance

Empirically

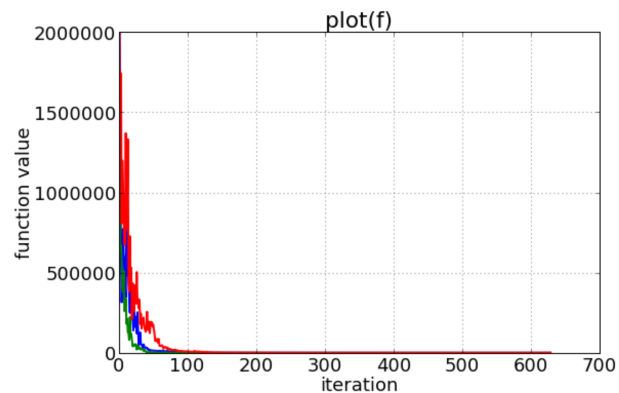
**convergence graphs** is all we have to start with

having the right presentation is important  
*too often neglected*

*the details are important*

62

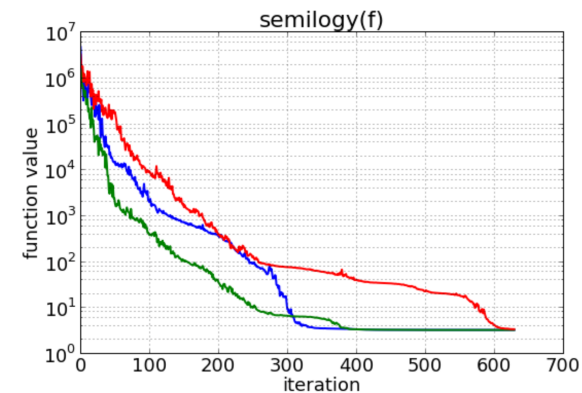
## Displaying Three Runs



not like this (it's unfortunately a common picture)

63

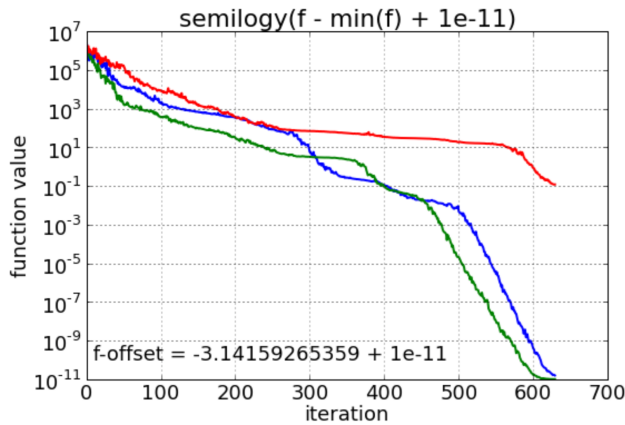
## Displaying Three Runs



better like this (shown are the same data),  
caveat: fails with negative f-values



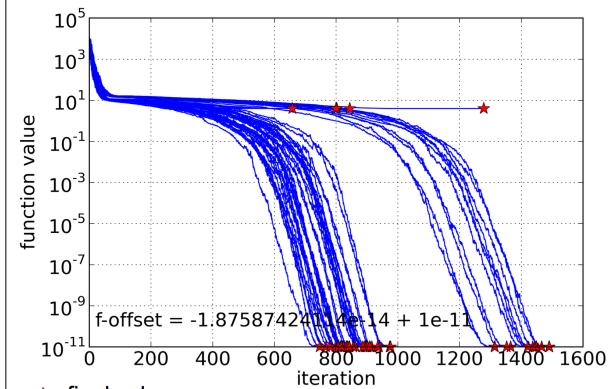
### Displaying Three Runs



even better like this: subtract minimum value over all runs

### Displaying 51 Runs

don't hesitate to display all data (the appendix is your friend)



\*: final value

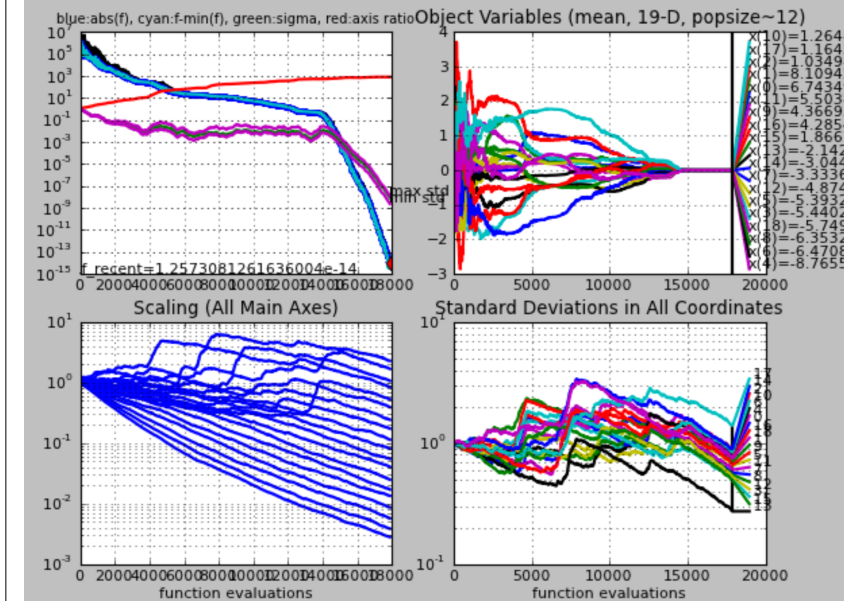
observation: three different "modes", which would be difficult to represent or recover in single statistics

### Visualization

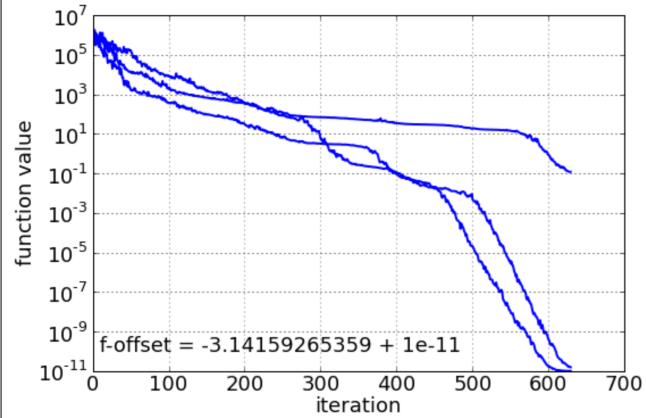
Other data than convergence graphs can be very instructive to visualize to understand performance

#### Example:

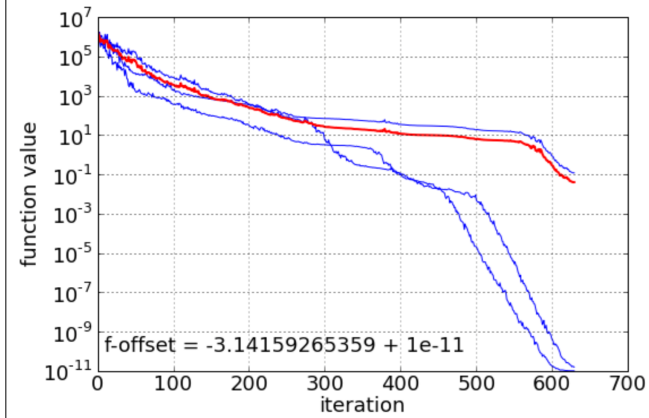
visualization of evolution of distribution in CMA-ES (see next slide)



### Which Statistics?



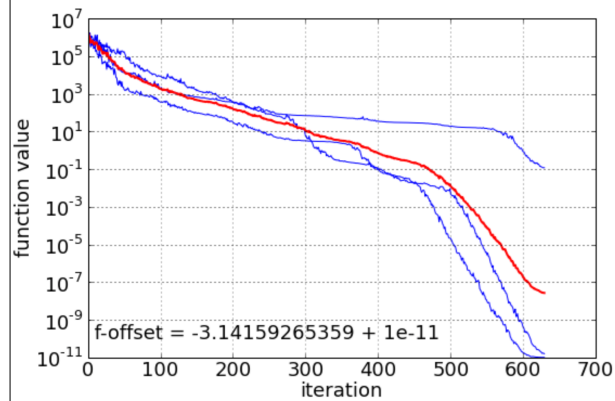
### Which Statistics?



mean/average function value

- tends to emphasize large values

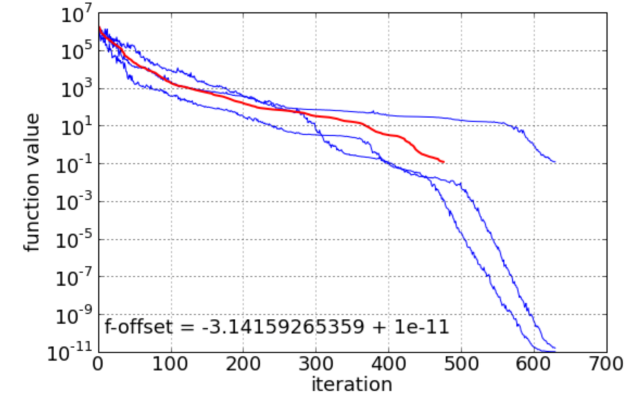
### Which Statistics?



geometric average function value  $\exp(\text{mean}_i(\log(f_i))) = (\prod_{i=1}^N f_i)^{1/N}$

- reflects "visual" average
- depends on offset
- artefact due to adding 1e-11

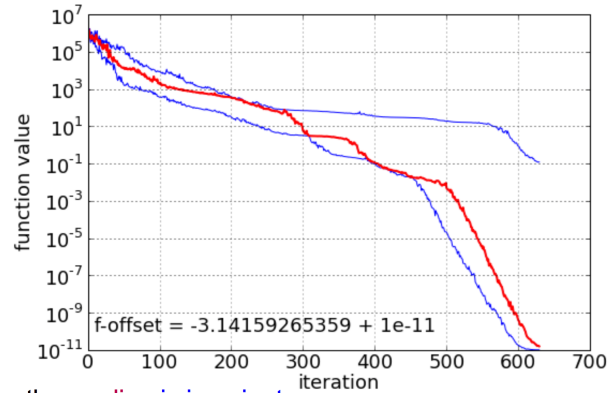
### Which Statistics?



average iterations

- reflects "visual" average
- here: incomplete

## Which Statistics?



the **median** is invariant

- unique for uneven number of data
- independent of log-scale, offset...  
 $\text{median}(\log(\text{data})) = \log(\text{median}(\text{data}))$
- same when taken over x- or y-direction

## Implication

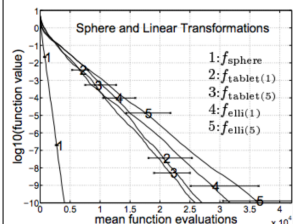
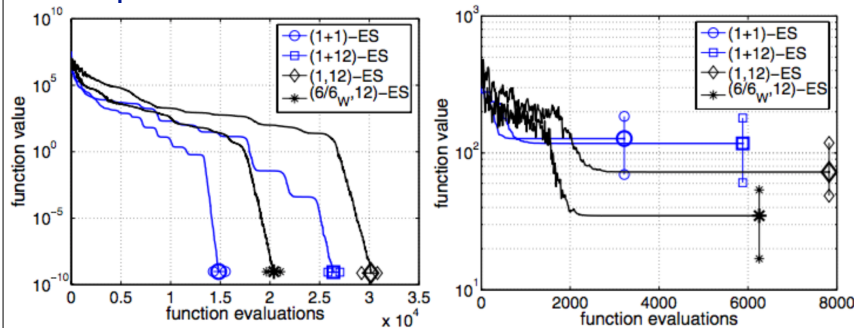
use the **median** as summary datum

more general: use quantiles as summary data

for example out of 15 data: 2nd, 8th, and 14th value represent the 10%, 50%, and 90%-tile

unless there are good reasons for a different statistics

## Examples



Comparison of 4 algorithms using the "median run" and the 90% central range of the final value on two different functions (Ellipsoid and Rastrigin)

caveat: this range display with simple error bars fails, if, e.g., 30% of all runs "converge"

## Statistical Assessment

1 Assess the meaning/**relevance** of a difference first (the only difficult part)

using enough data, any difference can be made significant

## Statistical Assessment

- ② Apply **rank-sum test** (Wilcoxon, Mann-Whitney U)  
only assumption: no equal data values

hypothesis:

compares  $sPr(x > y) \neq Pr(x < y) \neq 1/2$  ranking

two-sided 1%-significance p-value needs only 2x5 data values

For the same p-value, **fewer significant data** are better  
using enough data, any difference  
can be made significant

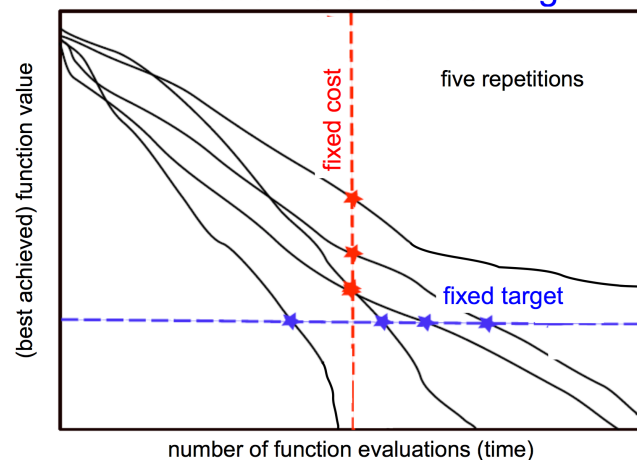
Generally: non-parametric tests, Kolmogorov-Smirnov test for ECDFs, no need to use the t-test

## Overview

- ① Problem Statement
  - Continuous Black-Box Optimization
  - Typical Difficulties
- ② Stochastic Black-Box Algorithms
  - General Template
  - Invariance
  - Comparisons of a few DFOs
- ③ Zoom on Evolution Strategies
  - Step-size Adaptation
  - Covariance Matrix Adaptation
- ④ Evaluating Black-Box Algorithms
  - Displaying results and visualization
  - Statistics
  - Average Runtime
  - Empirical Cumulative Distribution Function (ECDF)

## Measuring Performance from Convergence Graphs

fixed-cost versus fixed-target



## Evaluation of Search Algorithms

### Behind the scene

a performance should be

**quantitative** on the ratio scale (highest possible)

“algorithm A is two *times* better than algorithm B” is a meaningful statement

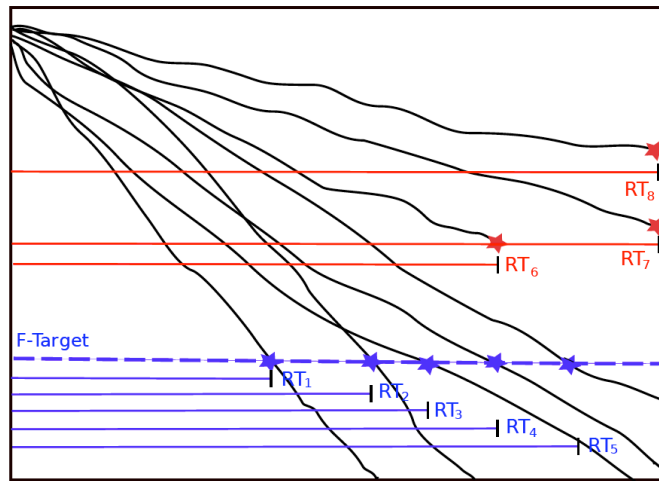
can assume a wide range of values

**meaningful (interpretable)** with regard to the real world

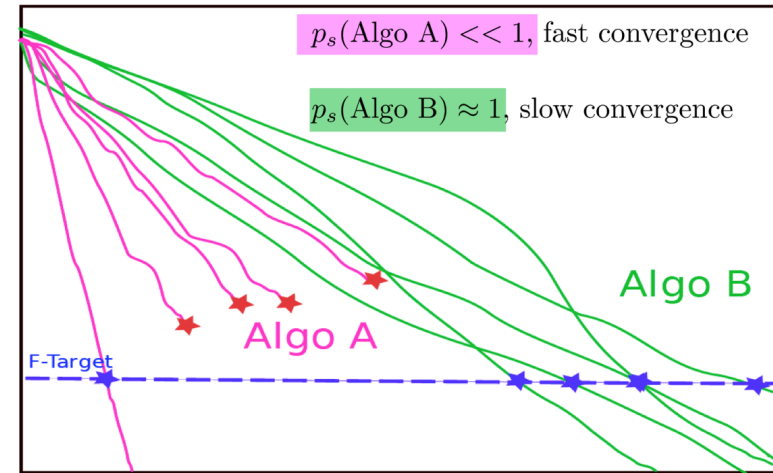
possible to transfer from benchmarking to real world

**runtime** or **first hitting time** is the prime candidate (we don't have many choices anyway)

### Collect Runtime to reach F-target



### Which Performance Measure to compare the two following scenario?



### Which Performance Measure

Algo Restart A:



$$p_s(\text{Algo Restart A}) = 1$$

Algo Restart B:



$$p_s(\text{Algo Restart B}) = 1$$

### Average Runtime (aRT)

$$\text{ERT} = \mathbb{E}[\text{RT}^r] = \frac{1-p_s}{p_s} \mathbb{E}[\text{RT}_{\text{unsuccessful}}] + \mathbb{E}[\text{RT}_{\text{successful}}]$$

aRT is an estimator for ERT

$$\text{aRT} = \frac{\#\text{Evals}}{\#\text{success}}$$

## Overview

### 1 Problem Statement

Continuous Black-Box Optimization  
Typical Difficulties

### 2 Stochastic Black-Box Algorithms

General Template  
Invariance  
Comparisons of a few DFOs

### 3 Zoom on Evolution Strategies

Step-size Adaptation  
Covariance Matrix Adaptation

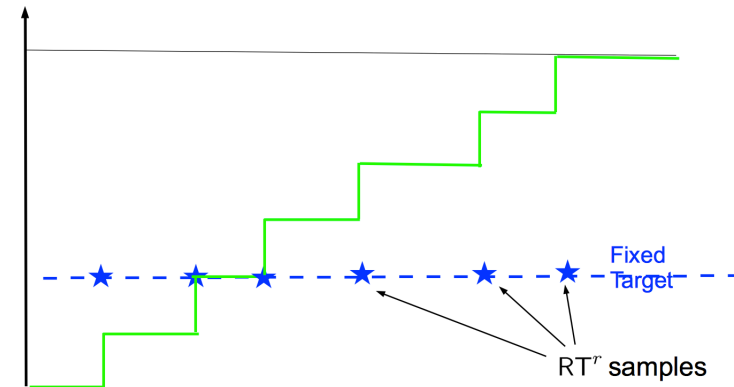
### 4 Evaluating Black-Box Algorithms

Displaying results and visualization  
Statistics  
Average Runtime  
Empirical Cumulative Distribution Function (ECDF)

85

## Empirical Cumulative Distribution Function

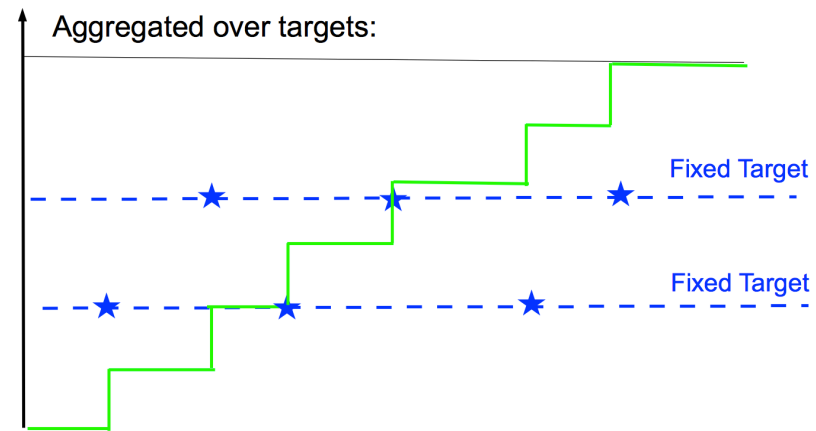
Single Target



86

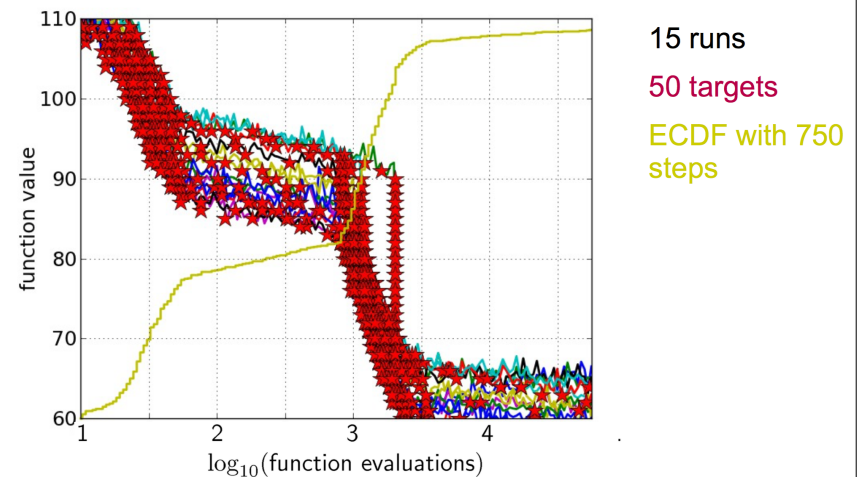
## Empirical Cumulative Distribution Function

Several Targets



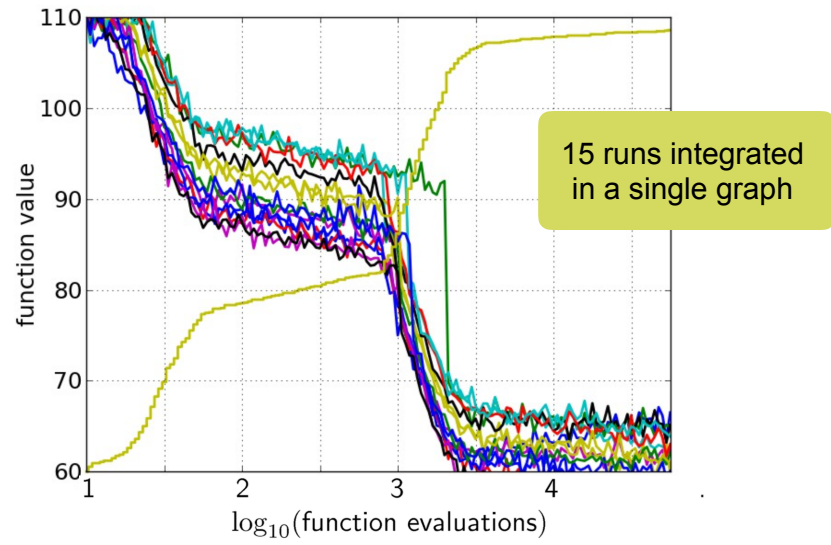
87

## Empirical Cumulative Distribution Function



88

### Empirical Cumulative Distribution Function



### Empirical Cumulative Distribution Function

