

Introduction to Randomized Continuous Optimization

Anne Auger & Nikolaus Hansen
Inria, Research Centre Saclay, France

anne.auger@inria.fr
nikolaus.hansen@lri.fr

<http://www.sigevo.org/gecco-2017/>

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).
GECCO '17 Companion, July 15-19, 2017, Berlin, Germany
© 2017 Copyright is held by the owner/author(s). ACM ISBN 978-1-4503-4939-0/17/07.
<http://dx.doi.org/10.1145/3067695.3067721>



1

Motivations and Objectives

Algorithms in **continuous domains** have **common grounds**
have to face the **same difficulties**
use **similar means** to overcome them

explicit or implicit variance control, ...

Teach you **basics** about **randomized optimization**
typical **difficulties**
important **algorithm design concepts**

avoid typical pitfalls

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

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

5

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

6

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

7

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

Navigation icons

8

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

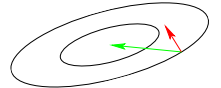
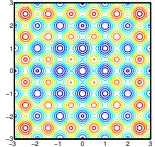
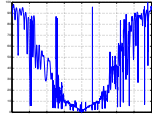
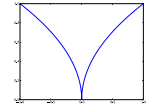
derivatives do not exist

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?

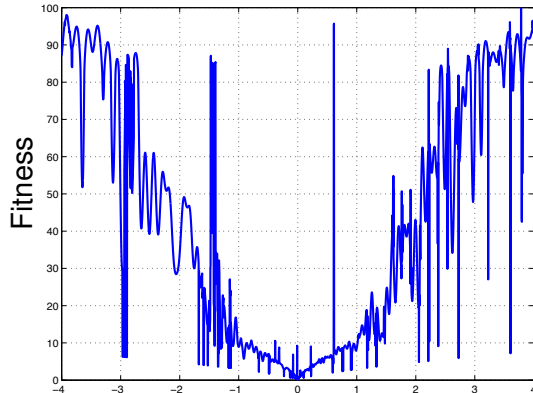
- non-linear, non-quadratic, non-convex
on linear and quadratic functions much better search policies are available
- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function
- dimensionality (size of search space)
(considerably) larger than three
- non-separability
dependencies between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

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.

13



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.

14



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.

15



Separable Problems

Definition (Separable Problem)

A function f is separable if

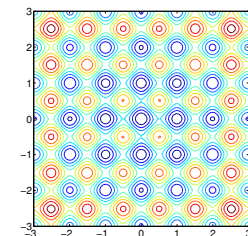
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left(\arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

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

Example: Additively decomposable functions

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

Rastrigin function



16



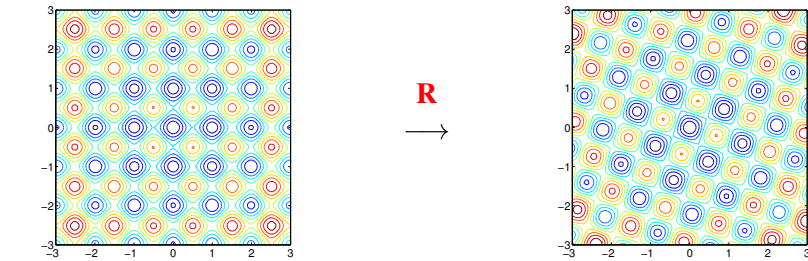
Non-Separable Problems

Building a non-separable problem from a separable one ^(1,2)

Rotating the coordinate system

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

R rotation matrix



¹ 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
² Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions: A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

What Makes a Function Difficult to Solve?

... 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	non-local policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed population-based method, stochastic, non-elitistic recombination operator serves as repair mechanism restarts

... metaphors

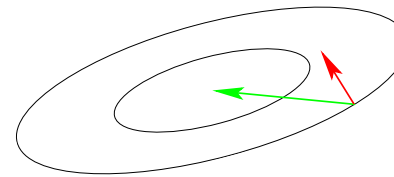
Ill-Conditioned Problems

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.

Metaphors

Evolutionary Computation		Optimization/Nonlinear Programmin
individual, offspring, parent	\longleftrightarrow	candidate solution decision variables design variables object variables
population	\longleftrightarrow	set of candidate solutions
fitness function	\longleftrightarrow	objective function loss function cost function error function
generation	\longleftrightarrow	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

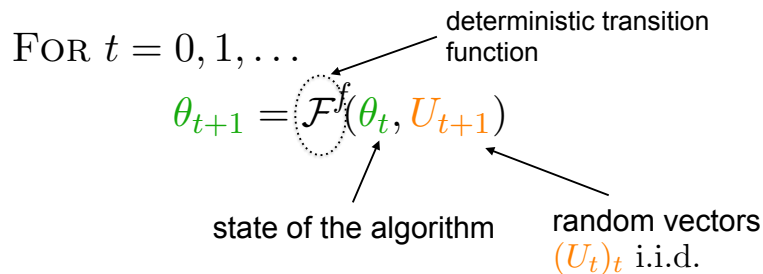
- Evolutionary Algorithms (continuous domain)
 - Differential Evolution [Storn & Price 1997]
 - Particle Swarm Optimization [Kennedy & Eberhart 1995]
 - Evolution Strategies, CMA-ES [Rechenberg 1965, Hansen & Ostermeier 2001]
 - Estimation of Distribution Algorithms (EDAs) [Larrañaga, Lozano, 2002]
 - 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

Stochastic / Randomized Algorithm

Iterative method



Optimization method

optimize $f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$
 θ_t typically encodes estimate(s) of the optimum of f

Example: Differential Evolution

[Storn, Price, 97]

$\theta_t = (X_t^1, \dots, X_t^N) \in (\mathbb{R}^n)^N$ population

Input CR $\in [0, 1]$, F $\in [0, 2]$, N pop size

For each $X_t \in \{X_t^1, \dots, X_t^N\}$

pick at random $X_t^{\alpha_1}, X_t^{\alpha_2}, X_t^{\alpha_3}$ (distinct from X_t)

sample $J = \text{Int}(1, \dots, n)$

for each coordinate $j = 1, \dots, n$

if $U_j(0, 1) < \text{CR}$ or $j = J$

$[Y]_j = [X_t^{\alpha_1}]_j + F ([X_t^{\alpha_2}]_j - [X_t^{\alpha_3}]_j)$

else

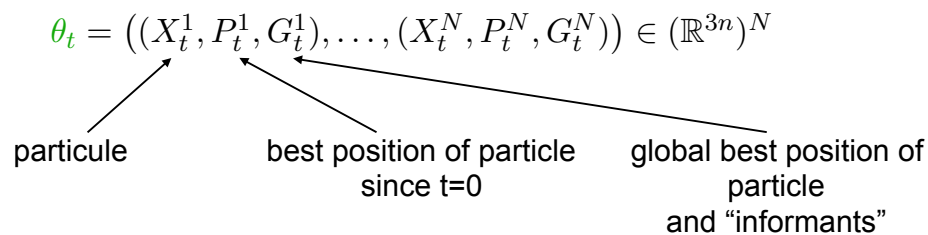
$[Y]_j = [X_t]_j$

if $f(Y) < f(X_t)$

$X_t \leftarrow Y$

ITERATION

Example 2: Particle Swarm Optimization



$$\theta_{t+1} = \mathcal{F}^f(\theta_t, \theta_{t-1}, U_{t+1})$$

25

Example 2: Particle Swarm Optimization (cont)

Input w inertia weight, N swarm size

For each particle X_t^k , for each coordinate j

$$[X_{t+\frac{1}{2}}^k]_j = [X_t^k]_j + [U_{t+1}^k]_j([P_t^k]_j - [X_t^k]_j) + [\tilde{U}_{t+1}^k]_j([G_t^k]_j - [X_t^k]_j)$$

$$[X_{t+1}^k]_j = [X_{t+\frac{1}{2}}^k]_j + w([X_t^k]_j - [X_{t-1}^k]_j)$$

$[U_{t+1}^k]_j, [\tilde{U}_{t+1}^k]_j$ random variables
i.i.d. $\sim \mathcal{U}(0, \varphi)$

ITERATION

26

Stochastic Search

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

Initialize distribution parameters θ , 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_λ 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_θ

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*

27

Evolution Strategies

New search points are sampled normally distributed

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

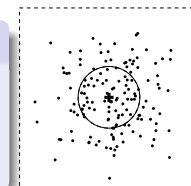
as perturbations of m , where $x_i, m \in \mathbb{R}^n, \sigma \in \mathbb{R}_+, \mathbf{C} \in \mathbb{R}^{n \times n}$

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** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

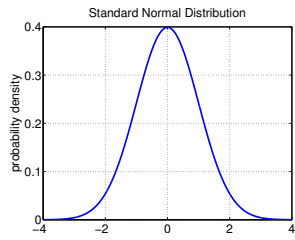
here, all new points are sampled with the same parameters

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

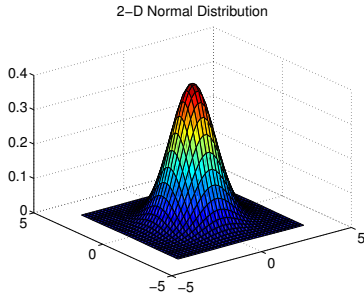


28

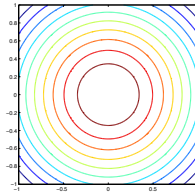
Normal Distribution



probability density of the 1-D standard normal distribution



probability density of a 2-D normal distribution

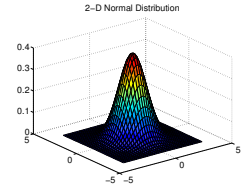


The Multi-Variate (n -Dimensional) Normal Distribution

Any multi-variate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ 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} .

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

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



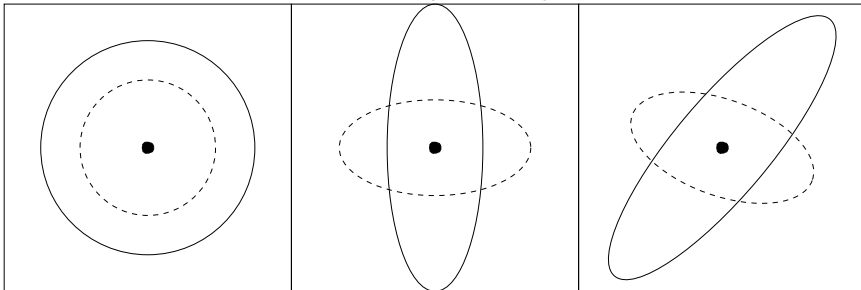
The **covariance matrix** \mathbf{C}

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



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

Lines of Equal Density



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

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

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

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



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

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

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

where

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

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

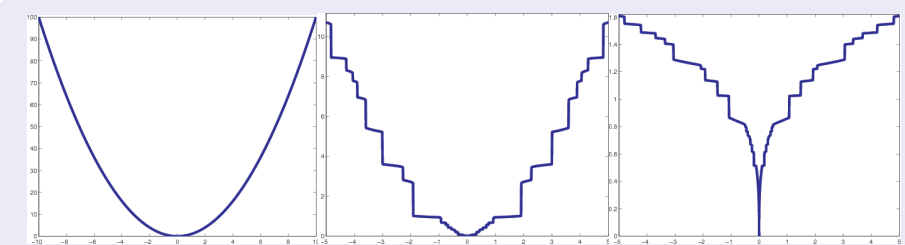


Invariance Under Monotonically Increasing Functions

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

3

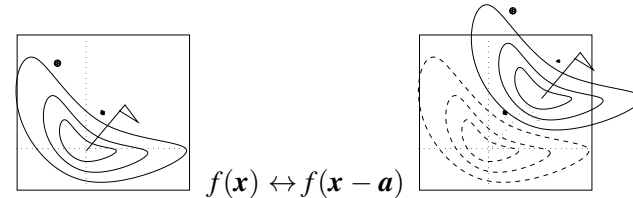
Whitley 1989. The GENITOR algorithm and selection pressure: Why rank-based allocation of reproductive trials is best, ICGA

33

Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



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

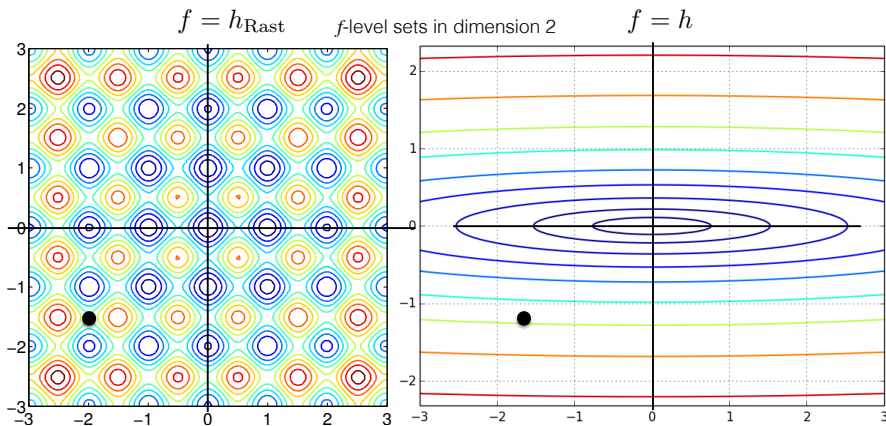
Identical behavior on f and f_a

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

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

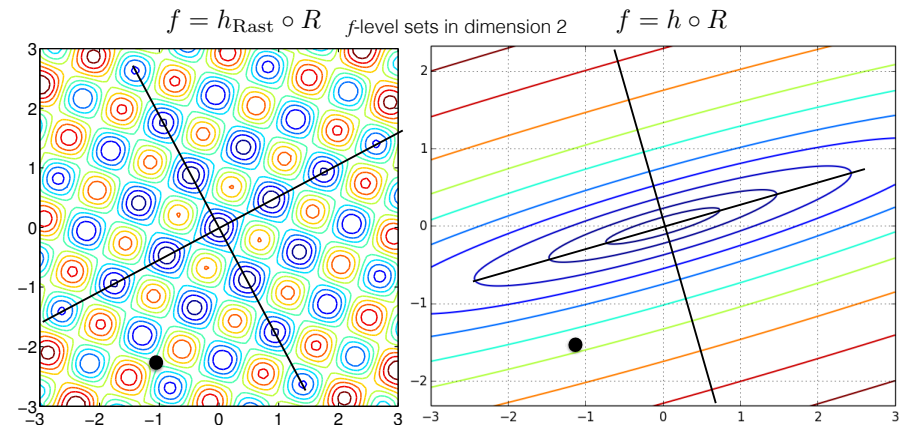
34

Invariance Under Rigid Search Space Transformations



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

Invariance Under Rigid Search Space Transformations



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

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

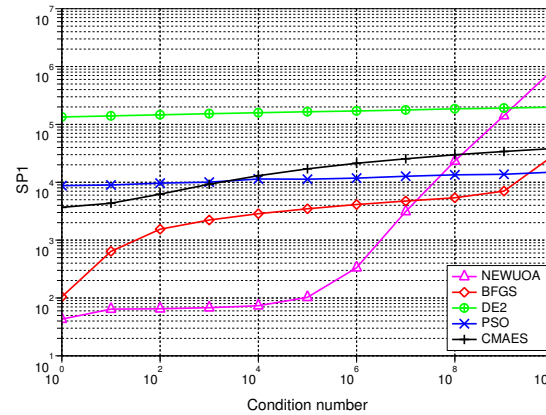
generalizing (identical) performance from a single function to a whole class of functions

consequently, invariance is important for the evaluation of search algorithms

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)

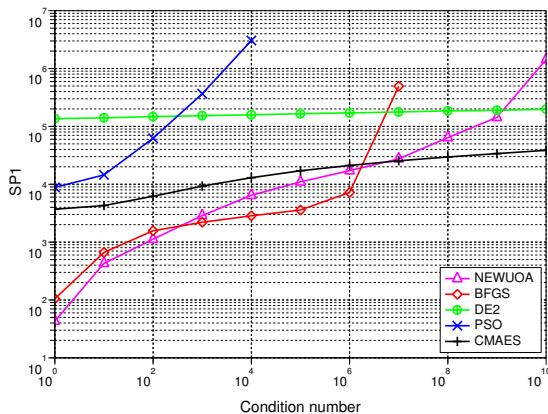
SP1 = average number of objective function evaluations¹⁴ to reach the target function value of $g^{-1}(10^{-9})$

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

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)

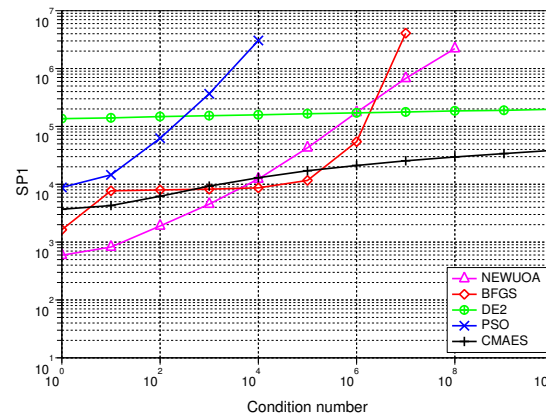
SP1 = average number of objective function evaluations¹⁵ to reach the target function value of $g^{-1}(10^{-9})$

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

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

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

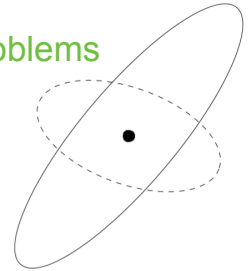
41

Zoom on ESs: Objectives

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

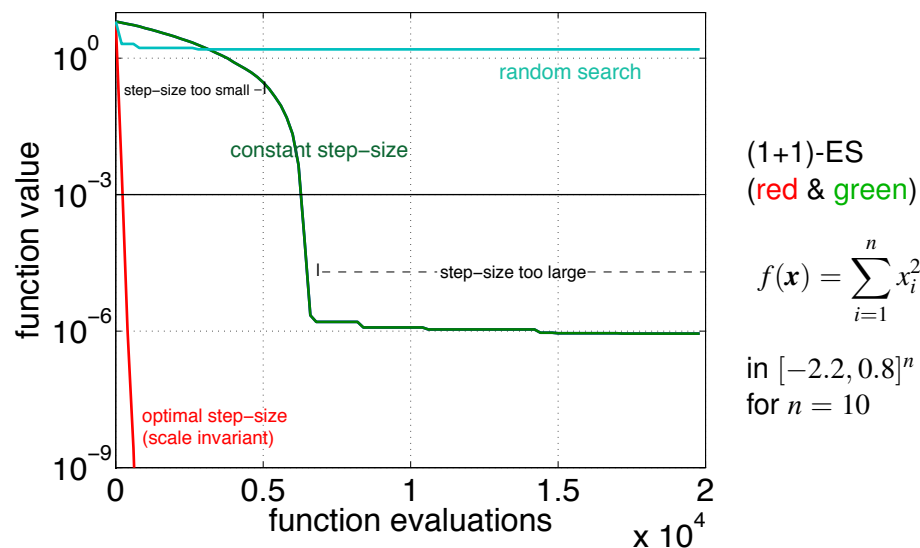
step-size control (overall standard deviation)
allows to achieve linear convergence

covariance matrix control
allows to solve ill-conditioned problems



42

Why Step-Size Control?



43

Methods for Step-Size Control

- **1/5-th success rule^{ab}**, often applied with “+”-selection
 - increase step-size if more than 20% of the new solutions are successful,
 - decrease otherwise
- **σ -self-adaptation^c**, 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^d** (Cumulative Step-size Adaptation, CSA)^e
 - self-adaptation derandomized and non-localized

^aRechenberg 1973, *Evolutionsstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog

^bSchumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

^cSchwefel 1981, *Numerical Optimization of Computer Models*, Wiley

^dHansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.*

9(2)

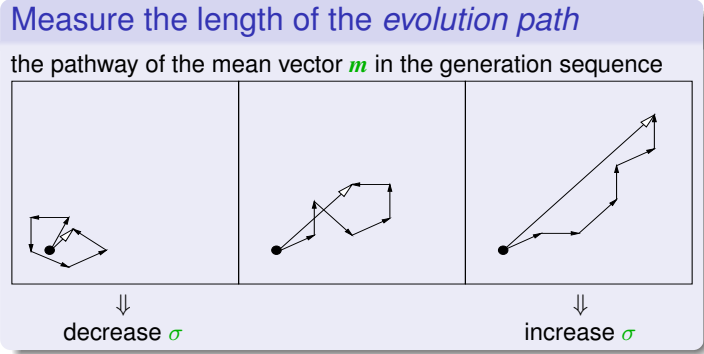
^eOstermeier *et al* 1994, Step-size adaptation based on non-local use of selection information, *PPSN IV*

44

Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

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



loosely speaking steps are

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

45



Path Length Control (CSA)

The Equations

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

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

46



Path Length Control (CSA)

The Equations

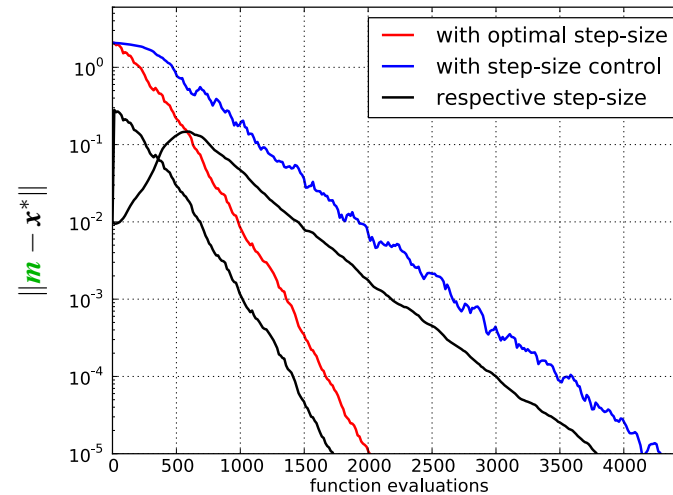
Initialize $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, evolution path $\mathbf{p}_\sigma = \mathbf{0}$,
set $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

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

47



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



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

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

48



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

49

Evolution Strategies

Recalling

New search points are sampled normally distributed

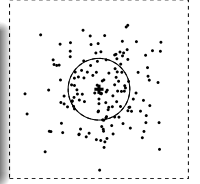
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

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} .



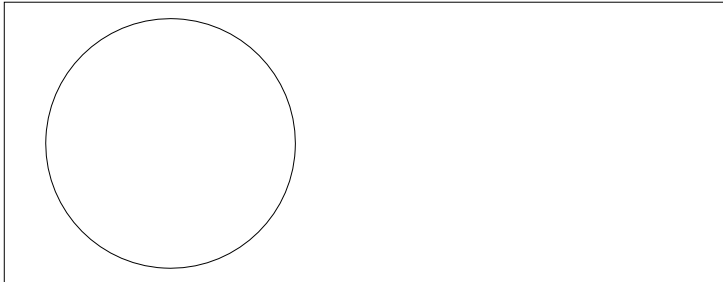
Navigation icons: back, forward, search, etc.

50

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

... equations

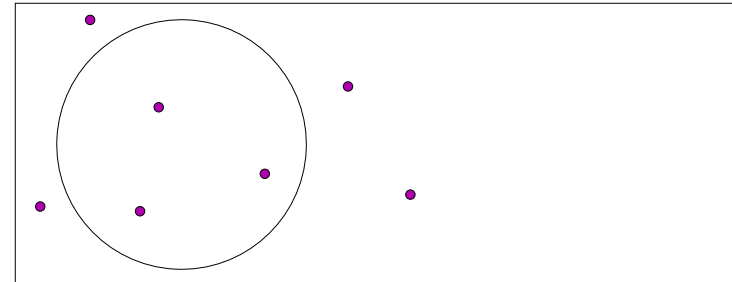
Navigation icons: back, forward, search, etc.

51

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

... equations

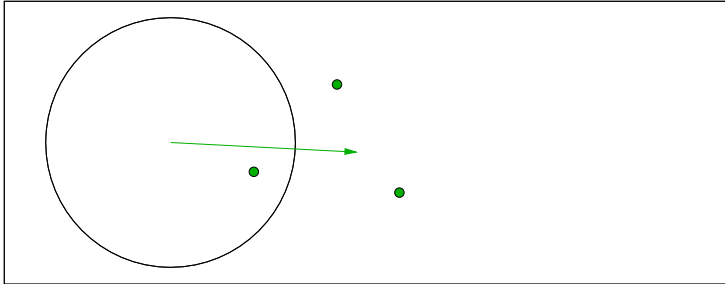
Navigation icons: back, forward, search, etc.

52

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)

... equations

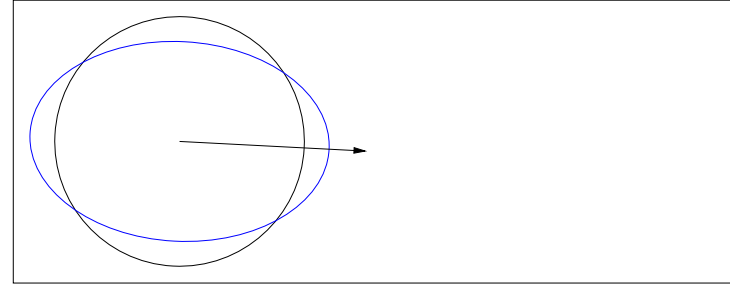


53

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

... equations

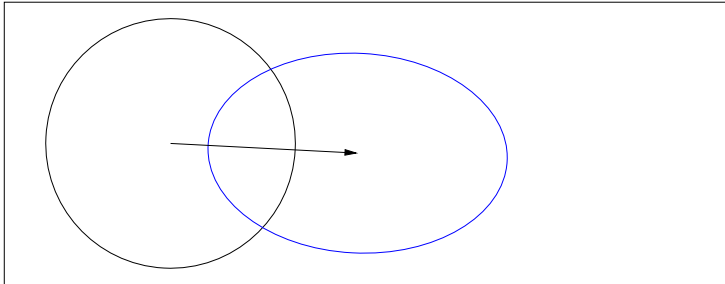


54

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)

... equations

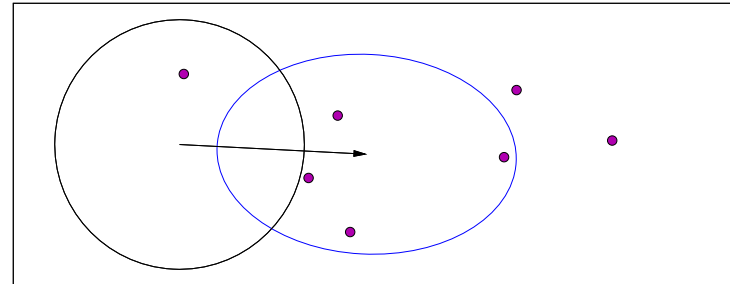


55

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)

... equations

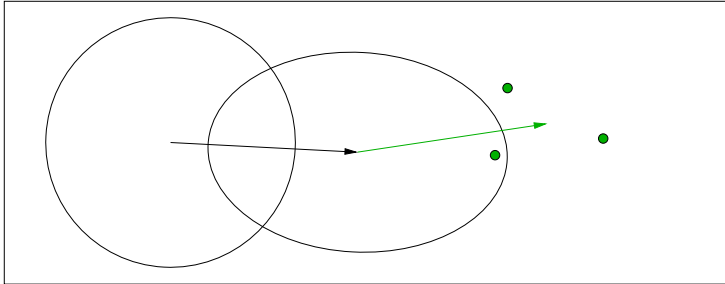


56

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}

... equations

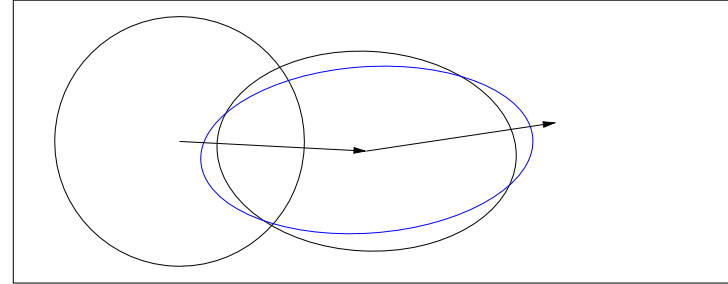


57

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

... equations

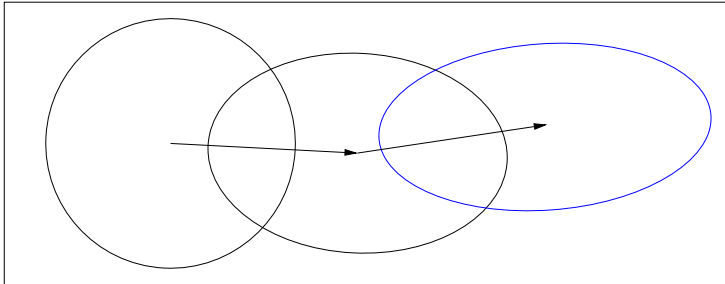


58

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



59

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

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

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

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

The rank-one update has been found independently in several domains^{6 7 8 9}

⁶ Kjellström & Taxén 1981. Stochastic Optimization in System Design, IEEE TCS

⁷ Hansen & Ostermeier 1996. Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation, ICEC

⁸ Ljung 1999. System Identification: Theory for the User

⁹ Haario et al 2001. An adaptive Metropolis algorithm, JSTOR

60



The CMA-ES

Input: $m \in \mathbb{R}^n, \sigma \in \mathbb{R}_+, \lambda$
Initialize: $C = I$, and $p_c = \mathbf{0}, p_\sigma = \mathbf{0}$,
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 $\mu_w = \frac{1}{\sum_{i=1}^\lambda w_i^2} \approx 0.3 \lambda$

While not terminate

$x_i = m + \sigma y_i, y_i \sim \mathcal{N}_i(\mathbf{0}, C)$, for $i = 1, \dots, \lambda$ sampling
 $m \leftarrow \sum_{i=1}^\mu w_i x_{i:\lambda} = m + \sigma y_w$ where $y_w = \sum_{i=1}^\mu w_i y_{i:\lambda}$ update mean
 $p_c \leftarrow (1 - c_c) p_c + \mathbf{1}_{\{\|p_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} y_w$ cumulation for C
 $p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} C^{-\frac{1}{2}} y_w$ cumulation for σ
 $C \leftarrow (1 - c_1 - c_\mu) C + c_1 p_c p_c^T + c_\mu \sum_{i=1}^\mu w_i y_{i:\lambda} y_{i:\lambda}^T$ update C
 $\sigma \leftarrow \sigma \times \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|p_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, I)\|} - 1\right)\right)$ update of σ

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

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

to the sphere model

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

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

without use of derivatives

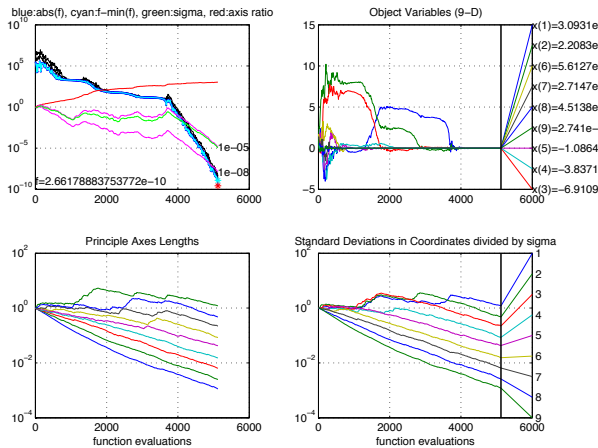
- lines of equal density align with lines of equal fitness

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

in a stochastic sense

Experimentum Crucis (1)

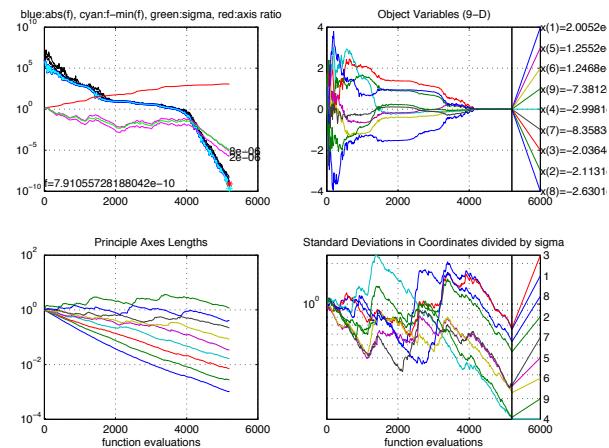
f convex quadratic, separable



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

Experimentum Crucis (2)

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



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

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