Continuous Optimization and CMA-ES

**Youhei Akimoto[1] & Anne Auger[2] & Nikolaus Hansen[2]**
**1. Shinshu University, Nagano, Japan**
**2. Inria Saclay Ile-de-France, France**

y_akimoto@shinshu-u.ac.jp
anne.auger@lri.fr
nikolaus.hansen@lri.fr

http://www.sigevo.org/gecco-2015/

1

We are happy to answer questions at any time.

2

# Content

3

# Problem Statement
Continuous Domain Search/Optimization

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

$$f : \mathcal{X} \subseteq \mathbb{R}^n \to \mathbb{R}, \qquad \boldsymbol{x} \mapsto f(\boldsymbol{x})$$

- 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
    . . . or to a robust solution $x$
  - ▶ solution $x$ with small function value $f(x)$ with least search cost
    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
    . . . or to a robust solution $x$
  - ▶ solution $x$ with small function value $f(x)$ with least search cost
    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
    . . . or to a robust solution $x$
  - ▶ solution $x$ with small function value $f(x)$ with least search cost
    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 $f : \mathcal{X} \subset \mathbb{R}^n \to \mathbb{R}$ to be *non-linear, non-separable* and to have at least moderate dimensionality, say $n \not\ll 10$.

Additionally, $f$ can be

- non-convex
- multimodal
  there are possibly many local optima
- 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

8

## Objective Function Properties

We assume $f : \mathcal{X} \subset \mathbb{R}^n \to \mathbb{R}$ to be *non-linear, non-separable* and to have at least moderate dimensionality, say $n \not\ll 10$.

Additionally, $f$ can be

- non-convex
- multimodal

there are possibly many local optima

- non-smooth

derivatives do not exist

- discontinuous, plateaus
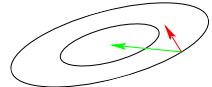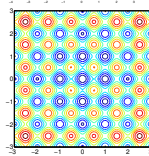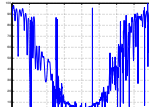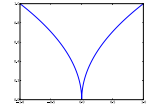- ill-conditioned
- noisy
- . . .

Goal: cope with any of these function properties

they are related to real-world problems

9

---

## 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

10

---

## Ruggedness

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



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

11

---

## 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.

12

315

## 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.

---

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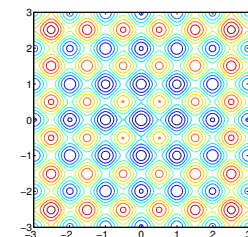
$$\arg \min_{(x_1, \ldots, x_n)} f(x_1, \ldots, x_n) = \left( \arg \min_{x_1} f(x_1, \ldots), \ldots, \arg \min_{x_n} f(\ldots, x_n) \right)$$

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

### Example: Additively decomposable functions

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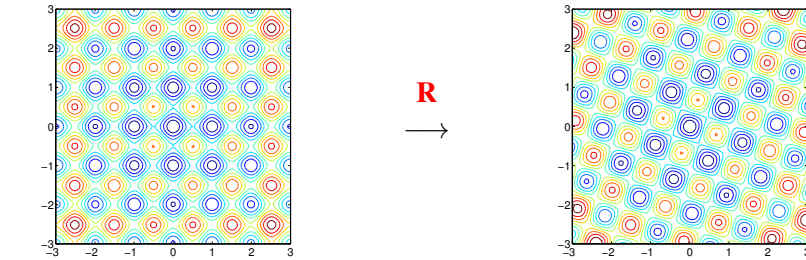
Rastrigin function

## Non-Separable Problems

Building a non-separable problem from a separable one [1,2]

### Rotating the coordinate system

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

$\mathbf{R}$ rotation matrix



$\mathbf{R}$

$\longrightarrow$

---

[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
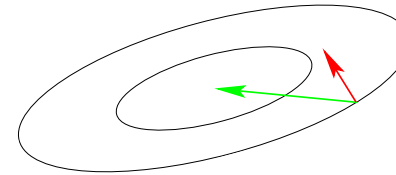
17

---

## Ill-Conditioned Problems

Curvature of level sets

Consider the convex-quadratic function
$$f(\boldsymbol{x}) = \tfrac{1}{2}(\boldsymbol{x}-\boldsymbol{x}^*)^T \boldsymbol{H}(\boldsymbol{x}-\boldsymbol{x}^*) = \tfrac{1}{2}\sum_i h_{i,i}\,(x_i-x_i^*)^2 + \tfrac{1}{2}\sum_{i\neq j} h_{i,j}\,(x_i-x_i^*)(x_j-x_j^*)$$

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



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

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

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

18

---

## What Makes a Function Difficult to Solve?

. . . and what can be done

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

. . . metaphors

19

---

## What Makes a Function Difficult to Solve?

. . . and what can be done

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

. . . metaphors

20

# What Makes a Function Difficult to Solve?
. . . and what can be done

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

. . . metaphors

21

---

Questions?

22

---

# Metaphors

| Evolutionary Computation | | Optimization/Nonlinear Programming |
|---|---|---|
| individual, offspring, parent | $\longleftrightarrow$ | candidate solution <br> decision variables <br> design variables <br> object variables |
| population | $\longleftrightarrow$ | set of candidate solutions |
| fitness function | $\longleftrightarrow$ | objective function <br> loss function <br> cost function <br> error function |
| generation | $\longleftrightarrow$ | iteration |

. . . methods: ESs

23

318

---

24

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta}, \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda, f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_\lambda))$

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*

25

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta}, \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda, f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_\lambda))$

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*

26

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta}, \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda, f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_\lambda))$

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*

27

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta}, \boldsymbol{x}_1, \ldots, \boldsymbol{x}_\lambda, f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_\lambda))$

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*

28

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta},\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda,f(\boldsymbol{x}_1),\ldots,f(\boldsymbol{x}_\lambda))$

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*

29

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta},\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda,f(\boldsymbol{x}_1),\ldots,f(\boldsymbol{x}_\lambda))$

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*

30

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta},\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda,f(\boldsymbol{x}_1),\ldots,f(\boldsymbol{x}_\lambda))$

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*

31

320

## Stochastic Search

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

Initialize distribution parameters $\boldsymbol{\theta}$, set population size $\lambda \in \mathbb{N}$
While not terminate

1. Sample distribution $P(\boldsymbol{x}|\boldsymbol{\theta}) \to \boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda \in \mathbb{R}^n$
2. Evaluate $\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda$ on $f$
3. Update parameters $\boldsymbol{\theta} \leftarrow F_\theta(\boldsymbol{\theta},\boldsymbol{x}_1,\ldots,\boldsymbol{x}_\lambda,f(\boldsymbol{x}_1),\ldots,f(\boldsymbol{x}_\lambda))$

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*

32

## The CMA-ES

Input: $m \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\lambda$
Initialize: $\mathbf{C} = \mathbf{I}$, 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 $w_{i=1\ldots\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^\mu w_i^2} \approx 0.3\,\lambda$

While not terminate

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

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

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

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

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

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

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

33

---

## Stochastic Search

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

Initialize distribution parameters $\theta$, set population size $\lambda \in \mathbb{N}$
While not terminate

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

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*

34

---

## Evolution Strategies

### New search points are sampled normally distributed

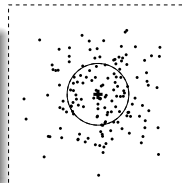$$x_i \sim m + \sigma\, \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \qquad \text{for } i = 1,\ldots,\lambda$$

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 $\sigma$.

35

---

## Evolution Strategies

### New search points are sampled normally distributed

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

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 $\sigma$.

36

## Why Normal Distributions?

1. widely observed in nature, for example as phenotypic traits

2. only stable distribution with finite variance

   stable means that the sum of normal variates is again normal:

   $$\mathcal{N}(\boldsymbol{x}, \mathbf{A}) + \mathcal{N}(\boldsymbol{y}, \mathbf{B}) \sim \mathcal{N}(\boldsymbol{x} + \boldsymbol{y}, \mathbf{A} + \mathbf{B})$$
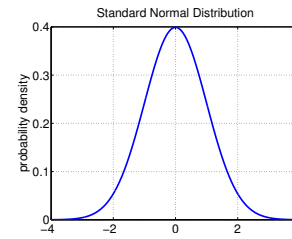
   helpful in design and analysis of algorithms related to the *central limit theorem*

3. most convenient way to generate isotropic search points

   the isotropic distribution does not favor any direction, rotational invariant

4. maximum entropy distribution with finite variance

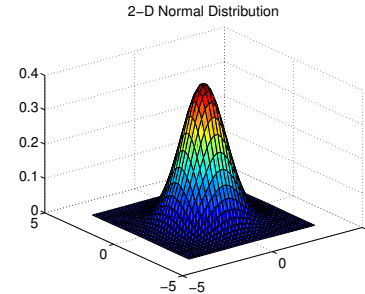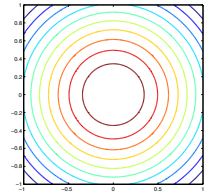   the least possible assumptions on $f$ in the distribution shape

---

## Normal Distribution



Standard Normal Distribution

probability density of the 1-D standard normal distribution



2–D Normal Distribution

probability density of a 2-D normal distribution

---

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

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

The mean value $\boldsymbol{m}$

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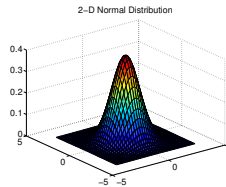


2–D Normal Distribution

---

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

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

The mean value $\boldsymbol{m}$

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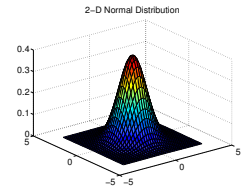- the distribution is symmetric about the distribution mean



2–D Normal Distribution

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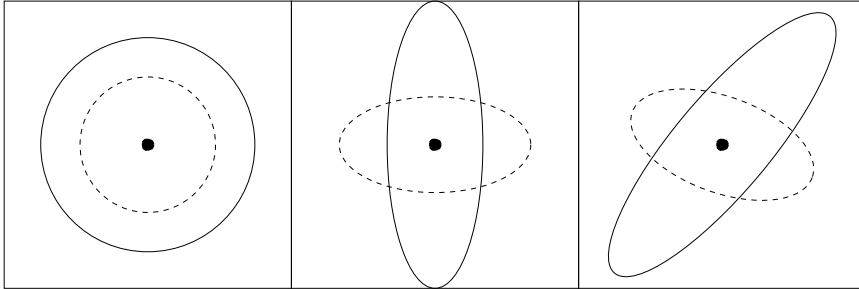
- determines the shape
- geometrical interpretation: any covariance matrix can be uniquely identified with the iso-density ellipsoid $\{\boldsymbol{x} \in \mathbb{R}^n \mid (\boldsymbol{x} - \boldsymbol{m})^{\mathrm{T}} \mathbf{C}^{-1} (\boldsymbol{x} - \boldsymbol{m}) = 1\}$

... 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



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

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

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

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

41

---

... 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



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

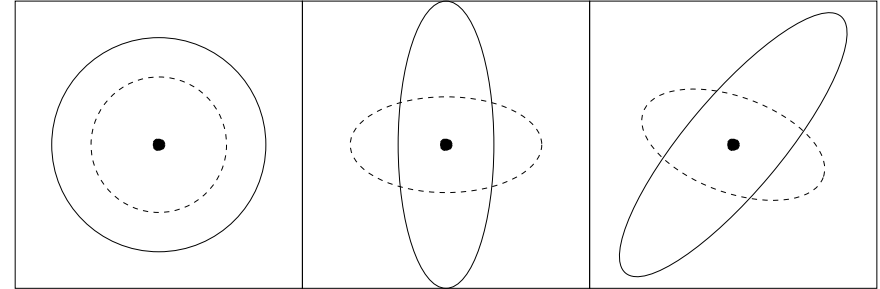$\mathcal{N}(m, D^2) \sim m + D\mathcal{N}(0, I)$
**n degrees of freedom**
components are
independent, scaled

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

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

42

---

... 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



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

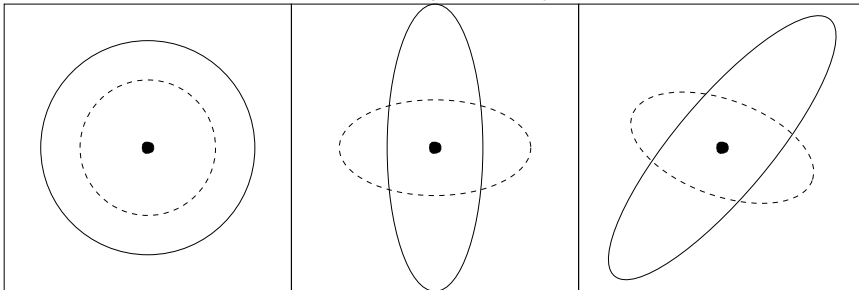$\mathcal{N}(m, D^2) \sim m + D\mathcal{N}(0, I)$
**n degrees of freedom**
components are
independent, scaled

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

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

43

---

# Effect of Dimensionality



$\|\mathcal{N}(0, I)\| \longrightarrow \mathcal{N}\left(\sqrt{n - 1/2}, 1/2\right)$ with modal value $\sqrt{n - 1}$

yet: maximum entropy distribution
also consider a difference between two vectors:
$\|\mathcal{N}(0, I) - \mathcal{N}(0, I)\| \sim \|\mathcal{N}(0, I) + \mathcal{N}(0, I)\| \sim \sqrt{2}\|\mathcal{N}(0, I)\|$

44

## Effect of Dimensionality



Norm of normally distributed vector

17D

2–D Normal Distribution

Norm of normally distributed vector

$\|\mathcal{N}(\mathbf{0}, \mathbf{I})\| \longrightarrow \mathcal{N}\left(\sqrt{n-1/2}, 1/2\right)$ with modal value $\sqrt{n-1}$

yet: maximum entropy distribution

also consider a difference between two vectors:
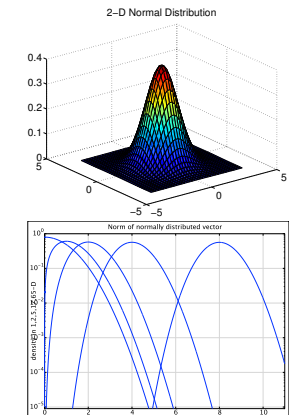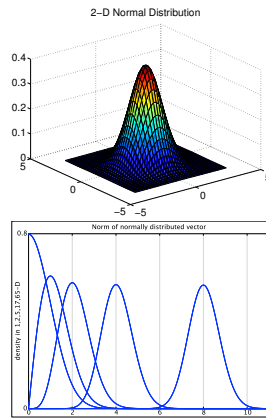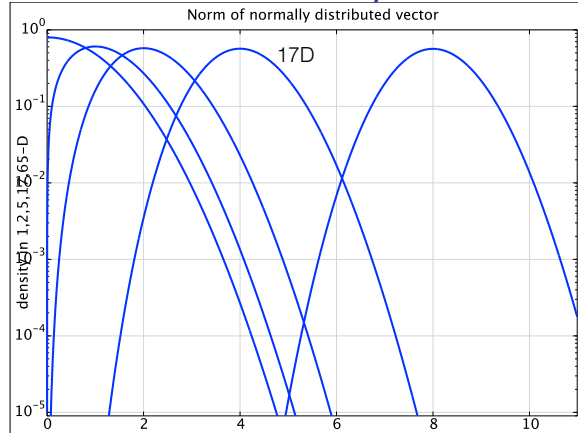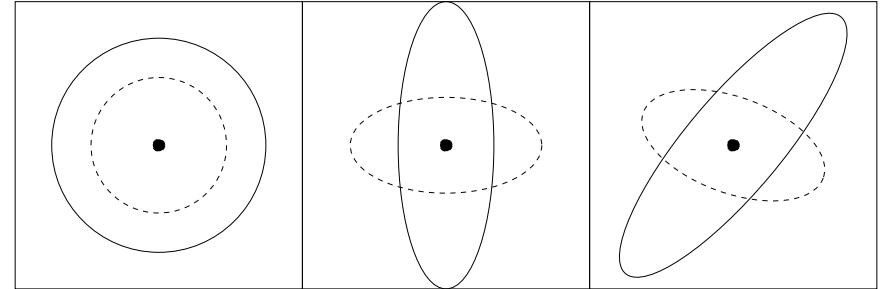
$\|\mathcal{N}(\mathbf{0}, \mathbf{I}) - \mathcal{N}(\mathbf{0}, \mathbf{I})\| \sim \|\mathcal{N}(\mathbf{0}, \mathbf{I}) + \mathcal{N}(\mathbf{0}, \mathbf{I})\| \sim \sqrt{2}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$

45

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

Lines of Equal Density



What is the implication for the distribution in this picture (considering large dimension)?

...ESs

46

## 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 $f(x_{1:\lambda}) \leq \cdots \leq f(x_{\lambda:\lambda})$.

The new mean reads

$$m \leftarrow \sum_{i=1}^{\mu} w_i x_{i:\lambda}$$

where

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

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

47

## 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 $f(x_{1:\lambda}) \leq \cdots \leq f(x_{\lambda:\lambda})$.

The new mean reads

$$m \leftarrow \sum_{i=1}^{\mu} w_i x_{i:\lambda}$$

where

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

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

48

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

Non-elitist selection and intermediate (weighted) recombination

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

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

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

where

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

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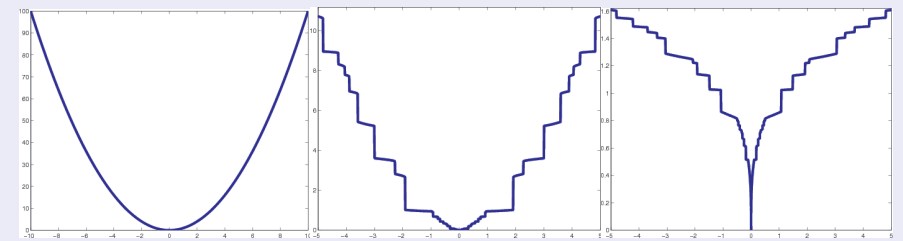and weighted intermediate recombination is applied.

49

## Invariance Under Monotonically Increasing Functions

### Rank-based algorithms

Update of all parameters uses only the ranks

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



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

$g$ is strictly monotonically increasing

$g$ preserves ranks

3

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

50

## Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



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

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

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

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

51

## Rotational Invariance in Search Space

- invariance to orthogonal (rigid) transformations $\mathbf{R}$, where $\mathbf{R}\mathbf{R}^{\mathrm{T}} = \mathbf{I}$

e.g. true for simple evolution strategies

recombination operators might jeopardize rotational invariance



$$f(\boldsymbol{x}) \leftrightarrow f(\mathbf{R}\boldsymbol{x})$$

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

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

45

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

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

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

52

## Landscape of Continuous Search Methods

*Gradient-based (Taylor, local)*
- Conjugate gradient methods [Fletcher & Reeves 1964]
- Quasi-Newton methods (BFGS) [Broyden et al 1970]

*Derivative-free optimization (DFO)*
- Trust-region methods (NEWUOA, BOBYQA) [Powell 2006, 2009]
- Simplex downhill [Nelder & Mead 1965]
- Pattern search [Hooke & Jeeves 1961, Audet & Dennis 2006]

*Stochastic (randomized) search methods*
- Evolutionary algorithms (broader sense, continuous domain)
  - Differential Evolution [Storn & Price 1997]
  - Particle Swarm Optimization [Kennedy & Eberhart 1995]
  - **Evolution Strategies** [Rechenberg 1965, Hansen & Ostermeier 2001]
- Simulated annealing [Kirkpatrick et al 1983]
- Simultaneous perturbation stochastic approximation (SPSA) [Spall 2000]

---

## 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

---

---

## Evolution Strategies
Recalling

New search points are sampled normally distributed

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

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 and $m \leftarrow \sum_{i=1}^{\mu} w_i \, x_{i:\lambda}$
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the *step length*
- the covariance matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the shape of the distribution ellipsoid

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

## Why Step-Size Control?



(1+1)-ES
(red & green)

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

in $[-2.2, 0.8]^n$
for $n = 10$

57

## Why Step-Size Control?

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



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

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

with optimal step-size $\sigma$

58

## Why Step-Size Control?

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



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

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

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

59

327

## Why Step-Size Control?

(5/5$_w$,10)-ES



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

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

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

60

## Why Step-Size Control?

(5/5$_w$,10)-ES



- with optimal step-size
- with step-size control
- respective step-size

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

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

comparing optimal versus default damping parameter $d_\sigma$: $\frac{1700}{1100} \approx 1.5$

61

## Why Step-Size Control?

$$\sigma_{\mathrm{opt}} = \sigma_{\mathrm{opt}}^* \frac{\|\boldsymbol{m}\|}{n} \approx \mu_w \frac{\|\boldsymbol{m}\|}{n}$$



$$-\frac{\varphi^*}{n}$$

$$\sigma_{\mathrm{opt}}^* \qquad \varphi^*$$

*evolution window* refers to the step-size interval (⊢——⊣) where reasonable performance is observed

62

## Methods for Step-Size Control

- 1/5-th success rule[a][b], often applied with "+"-selection

  increase step-size if more than $20\%$ of the new solutions are successful, decrease otherwise

- $\sigma$-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

---

[a] Rechenberg 1973, *Evolutionsstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog

[b] Schumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

[c] Schwefel 1981, *Numerical Optimization of Computer Models*, Wiley

[d] Hansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.* 9(2)

[e] Ostermeier *et al* 1994, Step-size adaptation based on non-local use of selection information, *PPSN IV*

63

328

## Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

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

### Measure the length of the *evolution path*

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



⇓          ⇓

decrease $\sigma$          increase $\sigma$

loosely speaking steps are

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

64

## Path Length Control (CSA)
The Equations

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

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

$$p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} y_w$$

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

---

## Path Length Control (CSA)
The Equations

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

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

$$p_\sigma \leftarrow (1 - c_\sigma) p_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} y_w$$

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

---

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



- with optimal step-size
- with step-size control
- respective step-size

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

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

---

1. Problem Statement

2. Evolution Strategies (ES)

3. Step-Size Control

4. Covariance Matrix Adaptation (CMA)
   - Covariance Matrix Rank-One Update
   - Cumulation—the Evolution Path
   - Covariance Matrix Rank-$\mu$ Update

5. CMA-ES Summary

6. Theoretical Foundations

7. Comparing Experiments

8. Summary and Final Remarks

# Evolution Strategies
Recalling

New search points are sampled normally distributed

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

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

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

# Covariance Matrix Adaptation
Rank-One Update

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

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

...equations

# Covariance Matrix Adaptation
Rank-One Update

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

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

...equations

# Covariance Matrix Adaptation
Rank-One Update

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

$y_w$, movement of the population mean $m$ (disregarding $\sigma$)

...equations

## Covariance Matrix Adaptation
Rank-One Update

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



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

## Covariance Matrix Adaptation
Rank-One Update

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



new distribution (disregarding $\sigma$)

## Covariance Matrix Adaptation
Rank-One Update

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



new distribution (disregarding $\sigma$)

## Covariance Matrix Adaptation
Rank-One Update

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



movement of the population mean $m$

## Covariance Matrix Adaptation
Rank-One Update

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



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

77

---

## Covariance Matrix Adaptation
Rank-One Update

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



new distribution,
$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times y_w y_w^{\mathrm{T}}$$
the ruling principle: the adaptation increases the likelihood of successful steps, $y_w$, to appear again
another viewpoint: the adaptation follows a natural gradient approximation of the expected fitness

78

---

## Covariance Matrix Adaptation
Rank-One Update

Initialize $m \in \mathbb{R}^n$, and $\mathbf{C} = \mathbf{I}$, set $\sigma = 1$, learning rate $c_{\mathrm{cov}} \approx 2/n^2$
While not terminate

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

$$m \;\leftarrow\; m + \sigma y_w \qquad \text{where } y_w = \sum_{i=1}^{\mu} w_i\, y_{i:\lambda}$$

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

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

79

---

$$\mathbf{C} \leftarrow (1 - c_{\mathrm{cov}})\mathbf{C} + c_{\mathrm{cov}}\mu_w y_w y_w^{\mathrm{T}}$$

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 $y_w$, 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

80

---

# 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

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

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.

---

# 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

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

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.

---

"Cumulation" is a widely used technique and also know as

- *exponential smoothing* in time series, forecasting
- exponentially weighted *mooving average*
- *iterate averaging* in stochastic approximation
- *momentum* in the back-propagation algorithm for ANNs
- . . .

"Cumulation" conducts a *low-pass* filtering, but there is more to it. . .

. . . why?

# Cumulation

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

Utilizing the Evolution Path

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



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

$$p_c \leftarrow \underbrace{(1 - c_c)\, p_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_c)^2}\sqrt{\mu_w}\, \mathbf{y}_w}_{\text{normalization factor}}$$

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + c_{\text{cov}} \underbrace{p_c p_c^{\mathrm{T}}}_{\text{rank-one}}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_c \ll 1$ such that $1/c_c$ is the "backward time horizon".

85

# Cumulation

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

Utilizing the Evolution Path

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



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

$$p_c \leftarrow \underbrace{(1 - c_c)\, p_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_c)^2}\sqrt{\mu_w}\, \mathbf{y}_w}_{\text{normalization factor}}$$

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + c_{\text{cov}} \underbrace{p_c p_c^{\mathrm{T}}}_{\text{rank-one}}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_c \ll 1$ such that $1/c_c$ is the "backward time horizon".

86

# Cumulation

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

Utilizing the Evolution Path

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



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

$$p_c \leftarrow \underbrace{(1 - c_c)\, p_c}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_c)^2}\sqrt{\mu_w}\, \mathbf{y}_w}_{\text{normalization factor}}$$

$$\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + c_{\text{cov}} \underbrace{p_c p_c^{\mathrm{T}}}_{\text{rank-one}}$$

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_{\text{cov}} \ll c_c \ll 1$ such that $1/c_c$ is the "backward time horizon".

87

Using an evolution path for the rank-one update of the covariance matrix reduces the number of function evaluations to adapt to a straight ridge from about $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.[a]

---

[a]Hansen & Auger 2013. Principled design of continuous stochastic search: From theory to practice.

Number of $f$-evaluations divided by dimension on the cigar function $f(\mathbf{x}) = x_1^2 + 10^6 \sum_{i=2}^{n} x_i^2$



$c_c = 1$ (no cumulation)

$c_c = 1/\sqrt{n}$

$c_c = 1/n$

The overall model complexity is $n^2$ but important parts of the model can be learned in time of order $n$

88

# Rank-$\mu$ Update

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

The rank-$\mu$ update extends the update rule for large population sizes $\lambda$ using $\mu > 1$ vectors to update $\mathbf{C}$ at each generation step.
The weighted empirical covariance matrix

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

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

with $\mu = \lambda$ weights can be negative [10]

The rank-$\mu$ update then reads

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

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

[10] Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.

89

---

# Rank-$\mu$ Update

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

The rank-$\mu$ update extends the update rule for large population sizes $\lambda$ using $\mu > 1$ vectors to update $\mathbf{C}$ at each generation step.
The weighted empirical covariance matrix

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

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

with $\mu = \lambda$ weights can be negative [10]

The rank-$\mu$ update then reads

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

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

[10] Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.

90

---

# Rank-$\mu$ Update

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

The rank-$\mu$ update extends the update rule for large population sizes $\lambda$ using $\mu > 1$ vectors to update $\mathbf{C}$ at each generation step.
The weighted empirical covariance matrix

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

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

with $\mu = \lambda$ weights can be negative [10]

The rank-$\mu$ update then reads

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

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

[10] Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.

91

335

---

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

$$\begin{aligned} \mathbf{C}_\mu &= \tfrac{1}{\mu}\sum \boldsymbol{y}_{i:\lambda}\boldsymbol{y}_{i:\lambda}^{\mathrm{T}} \\ \mathbf{C} &\leftarrow (1-1)\times\mathbf{C} + 1\times\mathbf{C}_\mu \end{aligned}$$

$m_{\mathrm{new}} \leftarrow \boldsymbol{m} + \tfrac{1}{\mu}\sum \boldsymbol{y}_{i:\lambda}$

new distribution

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

calculating $\mathbf{C}$ where $\mu = 50$, $w_1 = \cdots = w_\mu = \tfrac{1}{\mu}$, and $c_{\mathrm{cov}} = 1$

92

## Rank-$\mu$ CMA versus Estimation of Multivariate Normal Algorithm EMNA$_{\text{global}}$[11]



$x_i = m_{\text{old}} + y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$  $\quad \mathbf{C} \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - m_{\text{old}})(x_{i:\lambda} - m_{\text{old}})^{\mathrm{T}}$  $\quad m_{\text{new}} = m_{\text{old}} + \frac{1}{\mu} \sum y_{i:\lambda}$

rank-$\mu$ CMA conducts a PCA of steps

$x_i = m_{\text{old}} + y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$  $\quad \mathbf{C} \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - m_{\text{new}})(x_{i:\lambda} - m_{\text{new}})^{\mathrm{T}}$  $\quad m_{\text{new}} = m_{\text{old}} + \frac{1}{\mu} \sum y_{i:\lambda}$

EMNA$_{\text{global}}$ conducts a PCA of points

sampling of $\lambda = 150$ solutions (dots)

calculating $\mathbf{C}$ from $\mu = 50$ solutions

new distribution

$m_{\text{new}}$ is the minimizer for the variances when calculating $\mathbf{C}$

11 Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larranga, I. Inza and E. Bengoetxea (Eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102

---

The rank-$\mu$ update

- increases the possible learning rate in large populations
  roughly from $2/n^2$ to $\mu_w/n^2$

- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ [12]
  given $\mu_w \propto \lambda \propto n$

Therefore the rank-$\mu$ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3\, n + 10$

The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ .

Rank-one update and rank-$\mu$ update can be combined

... all equations

12 Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation, 11(1)*, pp. 1-18

---

The rank-$\mu$ update

- increases the possible learning rate in large populations
  roughly from $2/n^2$ to $\mu_w/n^2$

- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ [12]
  given $\mu_w \propto \lambda \propto n$

Therefore the rank-$\mu$ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3\, n + 10$

The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ .

Rank-one update and rank-$\mu$ update can be combined

... all equations

12 Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation, 11(1)*, pp. 1-18

---

The rank-$\mu$ update

- increases the possible learning rate in large populations
  roughly from $2/n^2$ to $\mu_w/n^2$

- can reduce the number of necessary generations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ [12]
  given $\mu_w \propto \lambda \propto n$

Therefore the rank-$\mu$ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3\, n + 10$

The rank-one update

- uses the evolution path and reduces the number of necessary function evaluations to learn straight ridges from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ .

Rank-one update and rank-$\mu$ update can be combined

... all equations

12 Hansen, Müller, and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation, 11(1)*, pp. 1-18

## Summary of Equations

The Covariance Matrix Adaptation Evolution Strategy

Input: $m \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\lambda$   (problem dependent)

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

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 $w_{i=1\ldots\lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^\mu w_i^2} \approx 0.3\,\lambda$

While not terminate

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

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

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

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

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

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

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

97

## Source Code Snippet

```
counteval = 0;  % the next 40 lines contain the 20 lines of interesting code
while counteval < stopeval

    % Generate and evaluate lambda offspring
    for k=1:lambda,
        arx(:,k) = xmean + sigma * B * (D .* randn(N,1)); % m + sig * Normal(0,C)
        arfitness(k) = feval(strfitnessfct, arx(:,k)); % objective function call
        counteval = counteval+1;
    end

    % Sort by fitness and compute weighted mean into xmean
    [arfitness, arindex] = sort(arfitness); % minimization
    xold = xmean;
    xmean = arx(:,arindex(1:mu))*weights;   % recombination, new mean value

    % Cumulation: Update evolution paths
    ps = (1-cs)*ps ...
        + sqrt(cs*(2-cs)*mueff) * invsqrtC * (xmean-xold) / sigma;
    hsig = norm(ps)/sqrt(1-(1-cs)^(2*counteval/lambda))/chiN < 1.4 + 2/(N+1);
    pc = (1-cc)*pc ...
        + hsig * sqrt(cc*(2-cc)*mueff) * (xmean-xold) / sigma;

    % Adapt covariance matrix C
    artmp = (1/sigma) * (arx(:,arindex(1:mu))-repmat(xold,1,mu));
    C = (1-c1-cmu) * C ...                  % regard old matrix
        + c1 * (pc*pc' ...                  % plus rank one update
            + (1-hsig) * cc*(2-cc) * C) ... % minor correction if hsig==0
        + cmu * artmp * diag(weights) * artmp'; % plus rank mu update

    % Adapt step size sigma
    sigma = sigma * exp((cs/damps)*(norm(ps)/chiN - 1));

    % Decomposition of C into B*diag(D.^2)*B' (diagonalization)
    if counteval - eigeneval > lambda/(c1+cmu)/N/10  % to achieve O(N^2)
        eigeneval = counteval;
        C = triu(C) + triu(C,1)';       % enforce symmetry
        [B,D] = eig(C);                 % eigen decomposition, B==normalized eigenvectors
        D = sqrt(diag(D));              % D is a vector of standard deviations now
        invsqrtC = B * diag(D.^-1) * B';
    end
```

98

## Strategy Internal Parameters

- related to selection and recombination
  - $\lambda$, offspring number, new solutions sampled, population size
  - $\mu$, parent number, solutions involved in updates of $m$, $\mathbf{C}$, and $\sigma$
  - $w_{i=1,\ldots,\mu}$, recombination weights
- related to $\mathbf{C}$-update
  - $c_c$, decay rate for the evolution path
  - $c_1$, learning rate for rank-one update of $\mathbf{C}$
  - $c_\mu$, learning rate for rank-$\mu$ update of $\mathbf{C}$
- related to $\sigma$-update
  - $c_\sigma$, decay rate of the evolution path
  - $d_\sigma$, damping for $\sigma$-change

Parameters were identified in carefully chosen experimental set ups. Parameters do not in the first place depend on the objective function and are not meant to be in the users choice.
Only(?) the population size $\lambda$ (and the initial $\sigma$) might be reasonably varied in a wide range, *depending on the objective function*
Useful: restarts with increasing population size (IPOP)

99

337

## Experimentum Crucis (0)

What did we want to achieve?

- reduce any convex-quadratic function

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

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

to the sphere model

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

without use of derivatives

- lines of equal density align with lines of equal fitness

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

in a stochastic sense
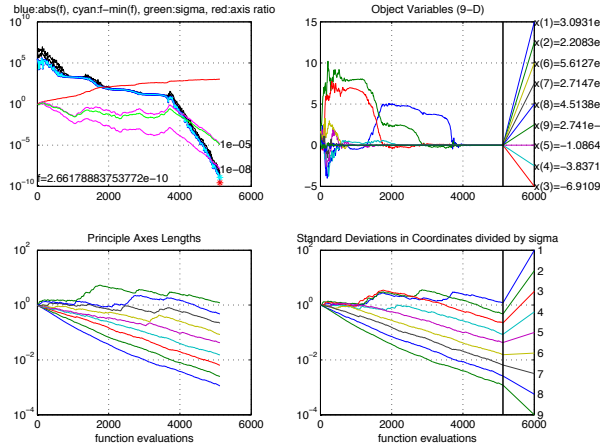
100

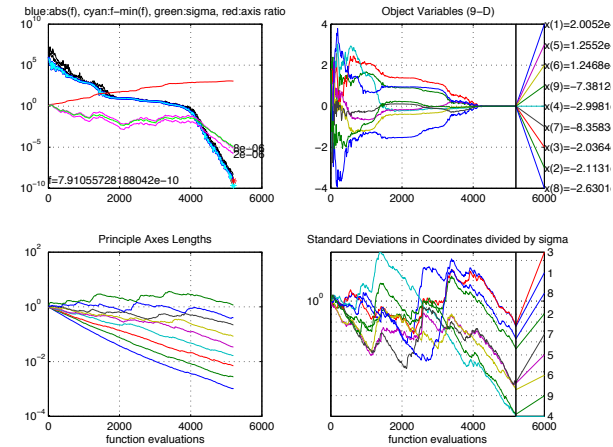## Experimentum Crucis (1)

*f* convex quadratic, separable



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

101

## Experimentum Crucis (2)

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



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

$$f(\boldsymbol{x}) = g\left(\boldsymbol{x}^{\mathrm{T}}\mathbf{H}\boldsymbol{x}\right), g : \mathbb{R} \to \mathbb{R} \text{ stricly increasing}$$

102

103

338

## Natural Gradient Descend

- Consider $\arg\min_{\boldsymbol{\theta}} \mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ under the sampling distribution $\boldsymbol{x} \sim p(.|\boldsymbol{\theta})$

  we could improve $\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ by following the gradient $\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$:

  $$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta}), \qquad \eta > 0$$

  $\nabla_{\boldsymbol{\theta}}$ depends on the parameterization of the distribution, therefore

- Consider the natural gradient of the expected transformed fitness

  $$\widetilde{\nabla}_{\boldsymbol{\theta}} \, \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta}) = F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta})$$

  $$= \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}} \ln p(\boldsymbol{x}|\boldsymbol{\theta}))$$

  using the Fisher information matrix $F_{\boldsymbol{\theta}} = \left(\left(E\frac{\partial^2 \log p(\boldsymbol{x}|\boldsymbol{\theta})}{\partial\theta_i\partial\theta_j}\right)\right)_{ij}$ of the density $p$.
  The natural gradient is invariant under re-parameterization of the distribution.

- A Monte-Carlo approximation reads

  $$\widetilde{\nabla}_{\boldsymbol{\theta}} \, \widehat{\mathrm{E}}(\widehat{w}(f(\boldsymbol{x}))|\boldsymbol{\theta}) = \sum_{i=1}^{\lambda} w_i F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}} \ln p(\boldsymbol{x}_{i:\lambda}|\boldsymbol{\theta}), \quad w_i = \widehat{w}(f(\boldsymbol{x}_{i:\lambda})|\boldsymbol{\theta})$$

104

# Natural Gradient Descend

- Consider $\arg\min_{\theta} \mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ under the sampling distribution $\boldsymbol{x} \sim p(.|\boldsymbol{\theta})$
  we could improve $\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ by following the gradient $\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta}), \qquad \eta > 0$$

  $\nabla_{\boldsymbol{\theta}}$ depends on the parameterization of the distribution, therefore

- Consider the natural gradient of the expected transformed fitness

$$\widetilde{\nabla}_{\boldsymbol{\theta}}\, \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta}) = F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta})$$
$$= \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}|\boldsymbol{\theta}))$$

  using the Fisher information matrix $F_{\boldsymbol{\theta}} = \left(\left(E\frac{\partial^2 \log p(x|\theta)}{\partial\theta_i\partial\theta_j}\right)\right)_{ij}$ of the density $p$.
  The natural gradient is invariant under re-parameterization of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\boldsymbol{\theta}}\,\widehat{\mathrm{E}}(\widehat{w}(f(\boldsymbol{x}))|\boldsymbol{\theta}) = \sum_{i=1}^{\lambda} w_i F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}_{i:\lambda}|\boldsymbol{\theta}), \quad w_i = \widehat{w}(f(\boldsymbol{x}_{i:\lambda})|\boldsymbol{\theta})$$

105

# Natural Gradient Descend

- Consider $\arg\min_{\theta} \mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ under the sampling distribution $\boldsymbol{x} \sim p(.|\boldsymbol{\theta})$
  we could improve $\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ by following the gradient $\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta}), \qquad \eta > 0$$

  $\nabla_{\boldsymbol{\theta}}$ depends on the parameterization of the distribution, therefore

- Consider the natural gradient of the expected transformed fitness

$$\widetilde{\nabla}_{\boldsymbol{\theta}}\, \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta}) = F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta})$$
$$= \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}|\boldsymbol{\theta}))$$

  using the Fisher information matrix $F_{\boldsymbol{\theta}} = \left(\left(E\frac{\partial^2 \log p(x|\theta)}{\partial\theta_i\partial\theta_j}\right)\right)_{ij}$ of the density $p$.
  The natural gradient is invariant under re-parameterization of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\boldsymbol{\theta}}\,\widehat{\mathrm{E}}(\widehat{w}(f(\boldsymbol{x}))|\boldsymbol{\theta}) = \sum_{i=1}^{\lambda} w_i F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}_{i:\lambda}|\boldsymbol{\theta}), \quad w_i = \widehat{w}(f(\boldsymbol{x}_{i:\lambda})|\boldsymbol{\theta})$$

106

# Natural Gradient Descend

- Consider $\arg\min_{\theta} \mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ under the sampling distribution $\boldsymbol{x} \sim p(.|\boldsymbol{\theta})$
  we could improve $\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ by following the gradient $\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta}), \qquad \eta > 0$$

  $\nabla_{\boldsymbol{\theta}}$ depends on the parameterization of the distribution, therefore

- Consider the natural gradient of the expected transformed fitness

$$\widetilde{\nabla}_{\boldsymbol{\theta}}\, \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta}) = F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta})$$
$$= \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}|\boldsymbol{\theta}))$$

  using the Fisher information matrix $F_{\boldsymbol{\theta}} = \left(\left(E\frac{\partial^2 \log p(x|\theta)}{\partial\theta_i\partial\theta_j}\right)\right)_{ij}$ of the density $p$.
  The natural gradient is invariant under re-parameterization of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\boldsymbol{\theta}}\,\widehat{\mathrm{E}}(\widehat{w}(f(\boldsymbol{x}))|\boldsymbol{\theta}) = \sum_{i=1}^{\lambda} w_i F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}_{i:\lambda}|\boldsymbol{\theta}), \quad w_i = \widehat{w}(f(\boldsymbol{x}_{i:\lambda})|\boldsymbol{\theta})$$

107

339

# Natural Gradient Descend

- Consider $\arg\min_{\theta} \mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ under the sampling distribution $\boldsymbol{x} \sim p(.|\boldsymbol{\theta})$
  we could improve $\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$ by following the gradient $\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta})$:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta\nabla_{\boldsymbol{\theta}}\mathrm{E}(f(\boldsymbol{x})|\boldsymbol{\theta}), \qquad \eta > 0$$

  $\nabla_{\boldsymbol{\theta}}$ depends on the parameterization of the distribution, therefore

- Consider the natural gradient of the expected transformed fitness

$$\widetilde{\nabla}_{\boldsymbol{\theta}}\, \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta}) = F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathrm{E}(w \circ P_f(f(\boldsymbol{x}))|\boldsymbol{\theta})$$
$$= \mathrm{E}(w \circ P_f(f(\boldsymbol{x}))F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}|\boldsymbol{\theta}))$$

  using the Fisher information matrix $F_{\boldsymbol{\theta}} = \left(\left(E\frac{\partial^2 \log p(x|\theta)}{\partial\theta_i\partial\theta_j}\right)\right)_{ij}$ of the density $p$.
  The natural gradient is invariant under re-parameterization of the distribution.

- A Monte-Carlo approximation reads

$$\tilde{\nabla}_{\boldsymbol{\theta}}\,\widehat{\mathrm{E}}(\widehat{w}(f(\boldsymbol{x}))|\boldsymbol{\theta}) = \sum_{i=1}^{\lambda} w_i F_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\ln p(\boldsymbol{x}_{i:\lambda}|\boldsymbol{\theta}), \quad w_i = \widehat{w}(f(\boldsymbol{x}_{i:\lambda})|\boldsymbol{\theta})$$

108

# CMA-ES = Natural Evolution Strategy + Cumulation

Natural gradient descend using the MC approximation and the normal distribution

- Rewriting the update of the distribution mean

$$m_{\text{new}} \leftarrow \sum_{i=1}^{\mu} w_i \, x_{i:\lambda} = m + \underbrace{\sum_{i=1}^{\mu} w_i (x_{i:\lambda} - m)}$$

$$\text{natural gradient for mean } \tilde{\tfrac{\partial}{\partial m}} \widehat{\mathrm{E}}(w \circ P_f(f(x)))|m, \mathbf{C})$$

- Rewriting the update of the covariance matrix[13]

$$\mathbf{C}_{\text{new}} \leftarrow \mathbf{C} + c_1 (\overbrace{p_c p_c^{\mathrm{T}}}^{\text{rank one}} - \mathbf{C})$$

$$+ \frac{c_\mu}{\sigma^2} \sum_{i=1}^{\mu} w_i \left( \overbrace{(x_{i:\lambda} - m)(x_{i:\lambda} - m)^{\mathrm{T}}}^{\text{rank-}\mu} - \sigma^2 \mathbf{C} \right)$$

$$\text{natural gradient for covariance matrix } \tilde{\tfrac{\partial}{\partial \mathbf{C}}} \widehat{\mathrm{E}}(w \circ P_f(f(x)))|m, \mathbf{C})$$

---

[13] Akimoto et al. (2010): Bidirectional Relation between CMA Evolution Strategies and Natural Evolution Strategies, PPSN XI

109

# Maximum Likelihood Update

**The new distribution mean $m$ maximizes the log-likelihood**

$$m_{\text{new}} = \arg\max_{m} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}}(x_{i:\lambda}|m)$$

independently of the given covariance matrix

$$\log p_{\mathcal{N}}(x|m, \mathbf{C}) = -\tfrac{1}{2} \log \det(2\pi \mathbf{C}) - \tfrac{1}{2}(x - m)^{\mathrm{T}} \mathbf{C}^{-1}(x - m)$$
$p_{\mathcal{N}}$ is the density of the multi-variate normal distribution

110

# Maximum Likelihood Update

**The new distribution mean $m$ maximizes the log-likelihood**

$$m_{\text{new}} = \arg\max_{m} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}}(x_{i:\lambda}|m)$$

independently of the given covariance matrix

**The rank-$\mu$ update matrix $\mathbf{C}_\mu$ maximizes the log-likelihood**

$$\mathbf{C}_\mu = \arg\max_{\mathbf{C}} \sum_{i=1}^{\mu} w_i \log p_{\mathcal{N}} \left( \frac{x_{i:\lambda} - m_{\text{old}}}{\sigma} \middle| m_{\text{old}}, \mathbf{C} \right)$$

$$\log p_{\mathcal{N}}(x|m, \mathbf{C}) = -\tfrac{1}{2} \log \det(2\pi \mathbf{C}) - \tfrac{1}{2}(x - m)^{\mathrm{T}} \mathbf{C}^{-1}(x - m)$$
$p_{\mathcal{N}}$ is the density of the multi-variate normal distribution

111

340

# Variable Metric

On the function class

$$f(x) = g\left( \frac{1}{2}(x - x^*)H(x - x^*)^{\mathrm{T}} \right)$$

the covariance matrix approximates the inverse Hessian up to a constant factor, that is:
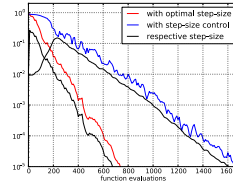
$$\mathbf{C} \propto H^{-1} \quad \text{(approximately)}$$

In effect, ellipsoidal level-sets are transformed into spherical level-sets.

$g : \mathbb{R} \to \mathbb{R}$ is strictly increasing

112

## On Convergence

Evolution Strategies converge with probability one on,
e.g., $g\left(\frac{1}{2}\boldsymbol{x}^{\mathrm{T}}\boldsymbol{H}\boldsymbol{x}\right)$ like

$$\|\boldsymbol{m}_k - \boldsymbol{x}^*\| \propto e^{-ck}, \qquad c \le \frac{0.25}{n}$$



Monte Carlo pure random search converges like

$$\|\boldsymbol{m}_k - \boldsymbol{x}^*\| \propto k^{-c} = e^{-c \log k}, \qquad c = \frac{1}{n}$$
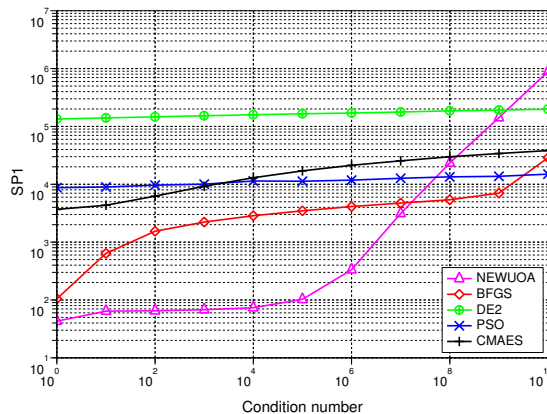
---

---

## 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)
NEWUAO (Powell 2004)
DE (Storn & Price 1996)
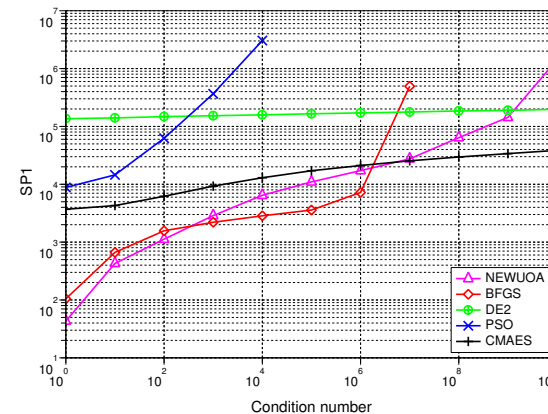PSO (Kennedy & Eberhart 1995)
CMA-ES (Hansen & Ostermeier 2001)

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

$\boldsymbol{H}$ diagonal

$g$ identity (for BFGS and NEWUOA)

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

SP1 = average number of objective function evaluations[14] 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 $\alpha$

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



BFGS (Broyden et al 1970)
NEWUAO (Powell 2004)
DE (Storn & Price 1996)
PSO (Kennedy & Eberhart 1995)
CMA-ES (Hansen & Ostermeier 2001)

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

$\boldsymbol{H}$ full

$g$ identity (for BFGS and NEWUOA)

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

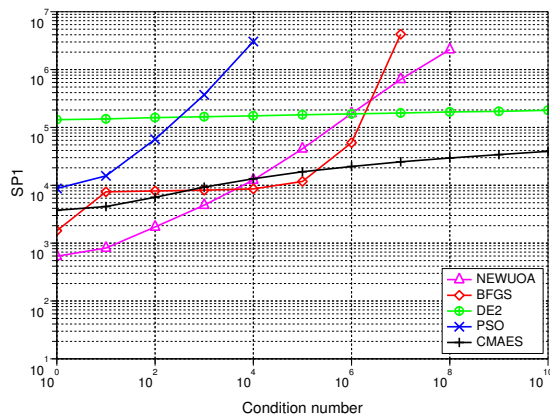SP1 = average number of objective function evaluations[15] to reach the target function value of $g^{-1}(10^{-9})$

[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 $\alpha$

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



BFGS (Broyden et al 1970)
NEWUAO (Powell 2004)
DE (Storn & Price 1996)
PSO (Kennedy & Eberhart 1995)
CMA-ES (Hansen & Ostermeier 2001)

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

$\boldsymbol{H}$ full

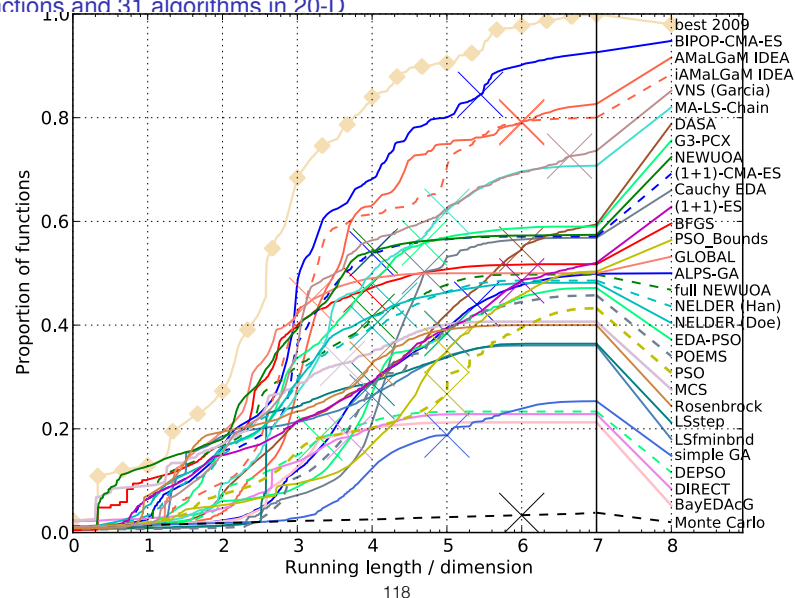$g : x \mapsto x^{1/4}$ (for BFGS and NEWUOA)

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

SP1 = average number of objective function evaluations[16] to reach the target function value of $g^{-1}(10^{-9})$

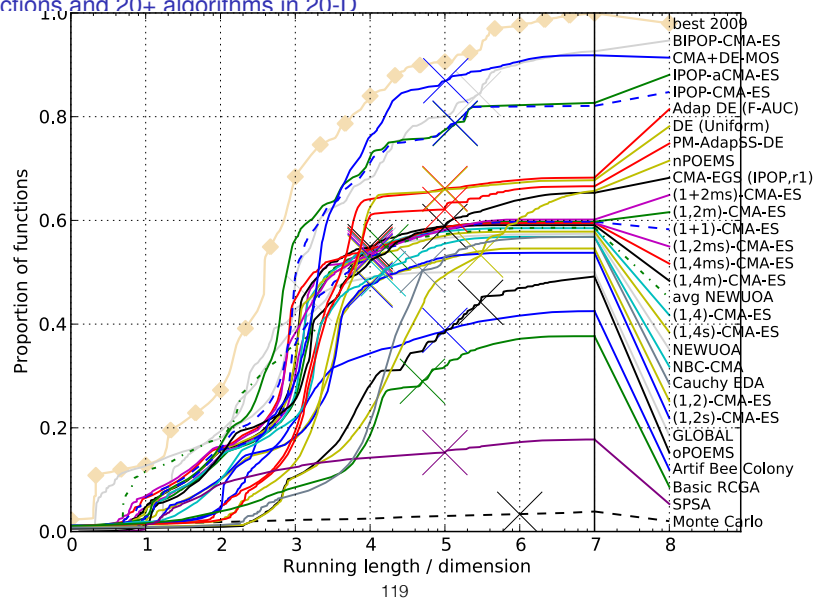[16] Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

117

# Comparison during BBOB at GECCO 2009
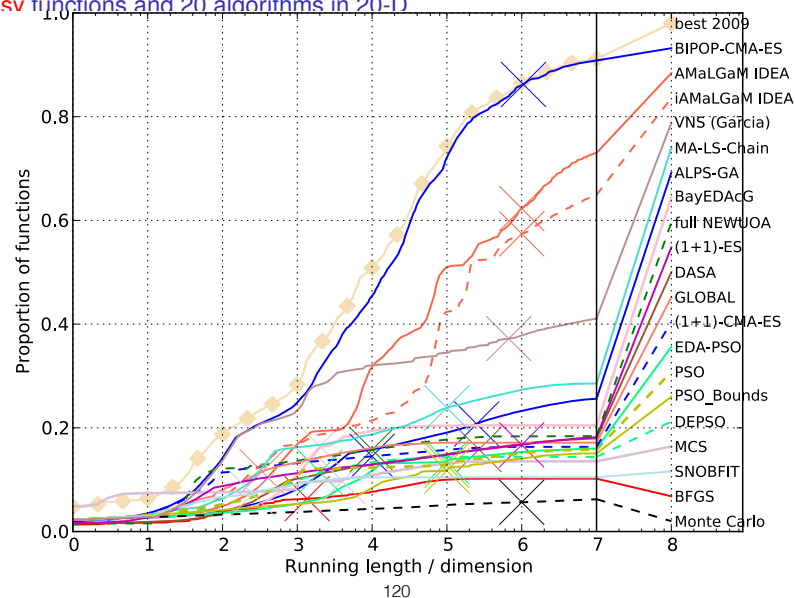
24 functions and 31 algorithms in 20-D



118

# Comparison during BBOB at GECCO 2010
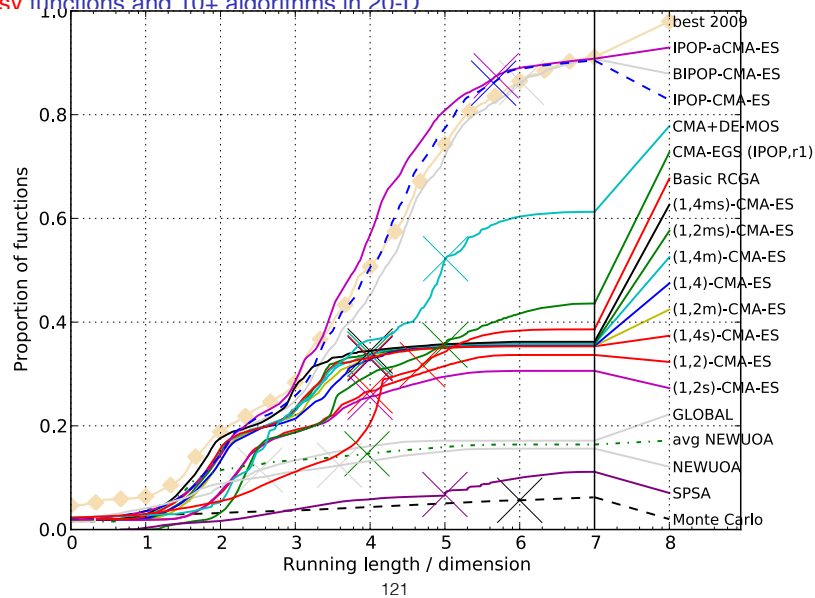
24 functions and 20+ algorithms in 20-D



119

# Comparison during BBOB at GECCO 2009

30 noisy functions and 20 algorithms in 20-D



120

342

## Comparison during BBOB at GECCO 2010

30 noisy functions and 10+ algorithms in 20-D



Legend (top to bottom):
- best 2009
- IPOP-aCMA-ES
- BIPOP-CMA-ES
- IPOP-CMA-ES
- CMA+DE-MOS
- CMA-EGS (IPOP,r1)
- Basic RCGA
- (1,4ms)-CMA-ES
- (1,2ms)-CMA-ES
- (1,4m)-CMA-ES
- (1,4)-CMA-ES
- (1,2m)-CMA-ES
- (1,4s)-CMA-ES
- (1,2)-CMA-ES
- (1,2s)-CMA-ES
- GLOBAL
- avg NEWUOA
- NEWUOA
- SPSA
- Monte Carlo

Axes: Proportion of functions (y), Running length / dimension (x)

---

---

## The Continuous Search Problem

Difficulties of a non-linear optimization problem are

- dimensionality and non-separabitity

  demands to exploit problem structure, e.g. neighborhood

  cave: design of benchmark functions

- ill-conditioning

  demands to acquire a second order model

- ruggedness

  demands a non-local (stochastic? population based?) approach

---

## 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 $g(f(\boldsymbol{x}))$ 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

   $\mathbf{C} \propto \boldsymbol{H}^{-1} \Longleftrightarrow$ adapts a variable metric

   $\Longleftrightarrow$ new (rotated) problem representation

   $\Longrightarrow f : \boldsymbol{x} \mapsto g(\boldsymbol{x}^{\mathrm{T}}\boldsymbol{H}\boldsymbol{x})$ reduces to $\boldsymbol{x} \mapsto \boldsymbol{x}^{\mathrm{T}}\boldsymbol{x}$

## 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

# Thank You

Source code for CMA-ES in C, Java, Matlab, Octave, Python, Scilab is available at `http://www.lri.fr/~hansen/cmaes_inmatlab.html`