## CMA-ES and Advanced Adaptation Mechanisms

**Youhei Akimoto[1]& Nikolaus Hansen[2]**
**1. Shinshu University, Nagano, Japan**
**2. Inria, Research Centre Saclay, France**

y_akimoto@shinshu-u.ac.jp
nikolaus.hansen@inria.fr

GECCO BERLIN 2017

1

---

We are happy to answer questions at any time.

2

---

## Topics

1. What makes the problem difficult to solve?

2. How does the CMA-ES work?

- Normal Distribution, Rank-Based Recombination
- Step-Size Adaptation (CSA)
- Covariance Matrix Adaptation (Hybrid-CMA)

3. What can/should the users do for the CMA-ES to work effectively on your problem?

- Restart, Increasing Population Size
- Restricted Covariance Matrix

3

---

## Topics

1. What makes the problem difficult to solve?

2. How does the CMA-ES work?

- Normal Distribution, Rank-Based Recombination
- Step-Size Adaptation (CSA)
- Covariance Matrix Adaptation (Hybrid-CMA)

3. What can/should the users do for the CMA-ES to work effectively on your problem?

- Restart, Increasing Population Size
- Restricted Covariance Matrix

4

# 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)

$$\text{x} \longrightarrow \boxed{\phantom{xxxx}} \longrightarrow \text{f(x)}$$

  - ▶ 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

---

# 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

---

# 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

---

# 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

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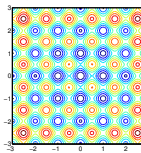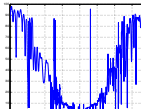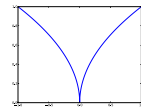
- discontinuous, plateaus
- ill-conditioned
- noisy
- . . .

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

7

---

---

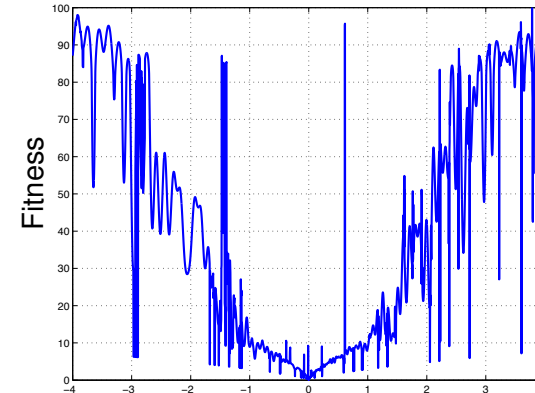# What Makes a Function Difficult to Solve?
Why stochastic search?

- 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

8

---

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

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

9

# 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.

10

---

## Separable Problems
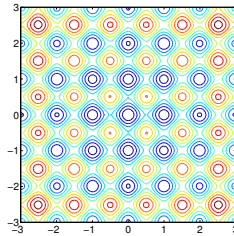
### Definition (Separable Problem)

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

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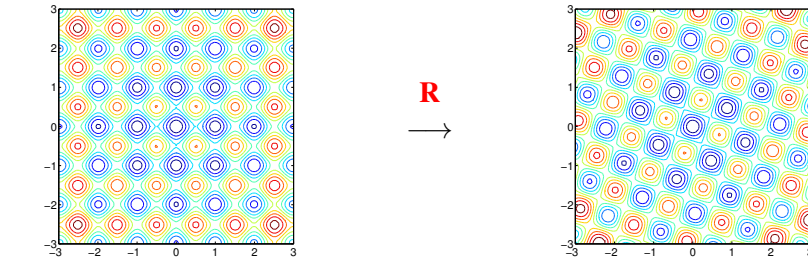
Rastrigin function

11

## 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

12

## 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 $\boldsymbol{H} \approx \boldsymbol{I}$ (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.

13

## 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

14

# 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> recombination operator <br> serves as repair mechanism <br><br> restarts |

... metaphors

14

---

# 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> recombination operator <br> serves as repair mechanism <br><br> restarts |

... metaphors

14

---

Questions?

15

---

# Topics

1. What makes the problem difficult to solve?

## 2. How does the CMA-ES work?

- Normal Distribution, Rank-Based Recombination
- Step-Size Adaptation (CSA)
- Covariance Matrix Adaptation (Hybrid-CMA)

3. What can/should the users do for the CMA-ES to work effectively on your problem?

- Restart, Increasing Population Size
- Restricted Covariance Matrix

16

## Stochastic Search

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

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

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

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*

17

---

## Stochastic Search

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

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

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

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*

18

---

## Stochastic Search

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

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

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

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*

19

---

## Stochastic Search

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

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

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

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*

20

## 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*

21

## 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*

22

## 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*

23

## 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*

24

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

While not terminate

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

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

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

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

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

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

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

25

---

## Evolution Strategies

New search points are sampled normally distributed

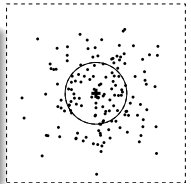$$x_i \sim m + \sigma\, \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \qquad \text{for } i = 1, \dots, \lambda$$

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 $\sigma$.

26

---

## Evolution Strategies

New search points are sampled normally distributed

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

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 $\sigma$.
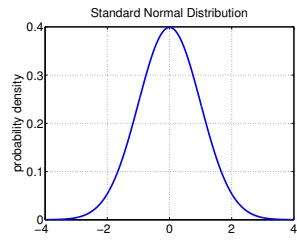
26

---

## 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:
   $$\mathcal{N}(x, \mathbf{A}) + \mathcal{N}(y, \mathbf{B}) \sim \mathcal{N}(x + y, \mathbf{A} + \mathbf{B})$$
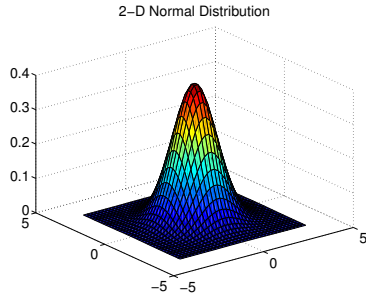   helpful in design and analysis of algorithms related to the *central limit theorem*
3. most convenient way to generate isotropic search points
   the isotropic distribution does not favor any direction, rotational invariant
4. maximum entropy distribution with finite variance
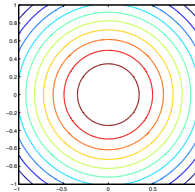   the least possible assumptions on $f$ in the distribution shape

27

# Normal Distribution



Standard Normal Distribution

probability density of the 1-D standard normal distribution



2–D Normal Distribution

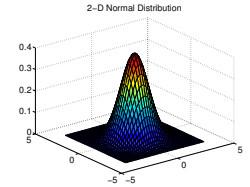probability density of a 2-D normal distribution



28

---

# 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 $\boldsymbol{m}$



2–D Normal Distribution

- 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 $\{\boldsymbol{x} \in \mathbb{R}^n \,|\, (\boldsymbol{x} - \boldsymbol{m})^{\mathrm{T}} \mathbf{C}^{-1} (\boldsymbol{x} - \boldsymbol{m}) = 1\}$

29

---

# 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 $\boldsymbol{m}$



2–D Normal Distribution

- 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 $\{\boldsymbol{x} \in \mathbb{R}^n \,|\, (\boldsymbol{x} - \boldsymbol{m})^{\mathrm{T}} \mathbf{C}^{-1} (\boldsymbol{x} - \boldsymbol{m}) = 1\}$

29

---

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

$$\begin{aligned}\mathcal{N}(\mathbf{0}, \mathbf{C}) &\sim \mathbf{A}\mathcal{N}(\mathbf{0}, \mathbf{I}) && \text{for any } \mathbf{A} \text{ s.t. } \mathbf{C} = \mathbf{A}\mathbf{A}^{\mathrm{T}} \\ &\sim \mathbf{B}\mathbf{D}\mathcal{N}(\mathbf{0}, \mathbf{I}) && \mathbf{C} = \mathbf{B}\mathbf{D}^2\mathbf{B}^{\mathrm{T}} \text{ (Eigen decomposition of } \mathbf{C}) \\ &\sim \mathcal{N}_1(0,1)d_1\boldsymbol{b}_1 + \cdots + \mathcal{N}_n(0,1)d_n\boldsymbol{b}_n \end{aligned}$$

$d_i$: square root of the eigenvalue of $\mathbf{C}$

$\boldsymbol{b}_i$: eigenvector of $\mathbf{C}$, corresponding to $d_i$


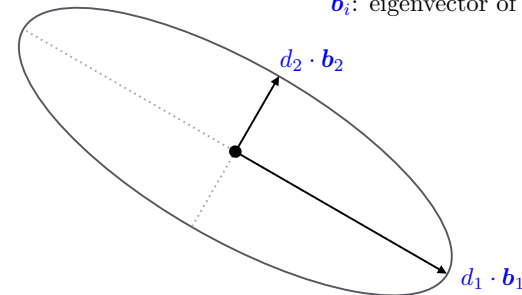
$d_2 \cdot \boldsymbol{b}_2$

$d_1 \cdot \boldsymbol{b}_1$

30

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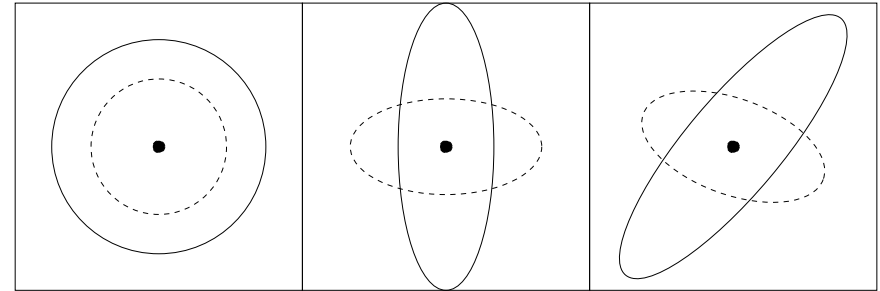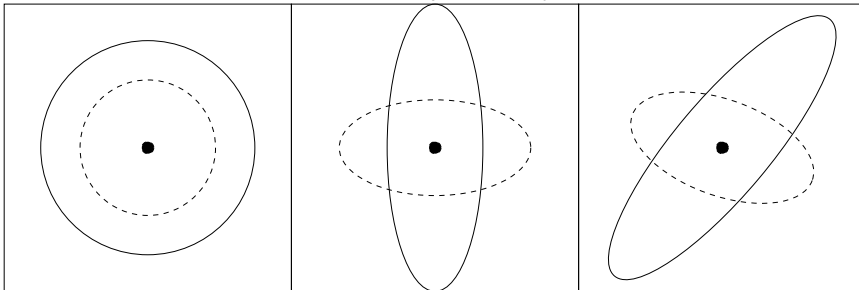
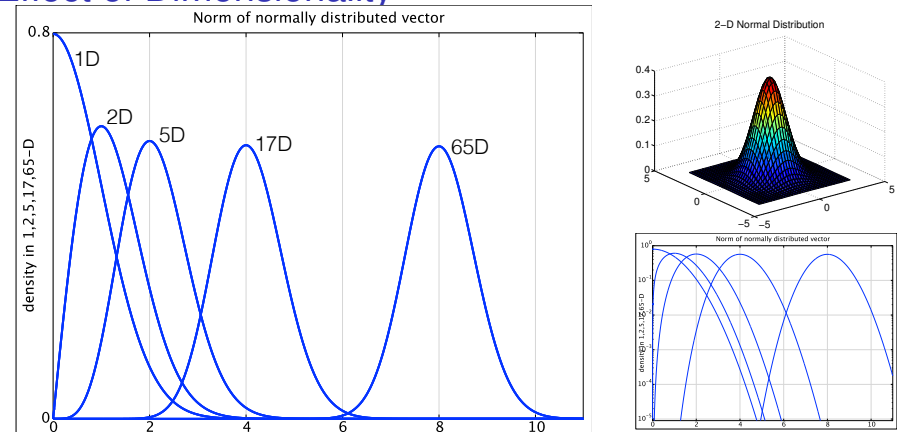### Lines of Equal Density



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

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

$\mathcal{N}(m, \mathbf{C}) \sim m + \mathbf{C}^{\frac{1}{2}}\mathcal{N}(\mathbf{0}, \mathbf{I})$
$(n^2 + n)/2$ 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 $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^\mathrm{T})$ holds for all $\mathbf{A}$.

31

---

---

---

## Effect of Dimensionality



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

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

32

# Effect of Dimensionality

Norm of normally distributed vector



17D

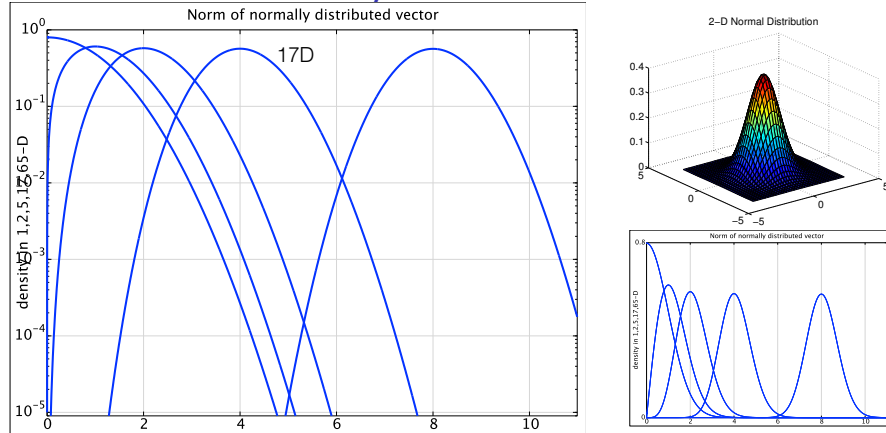2–D Normal Distribution

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:

$\|\mathcal{N}(\mathbf{0}, \mathbf{I}) - \mathcal{N}(\mathbf{0}, \mathbf{I})\| \sim \|\mathcal{N}(\mathbf{0}, \mathbf{I}) + \mathcal{N}(\mathbf{0}, \mathbf{I})\| \sim \sqrt{2}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|$
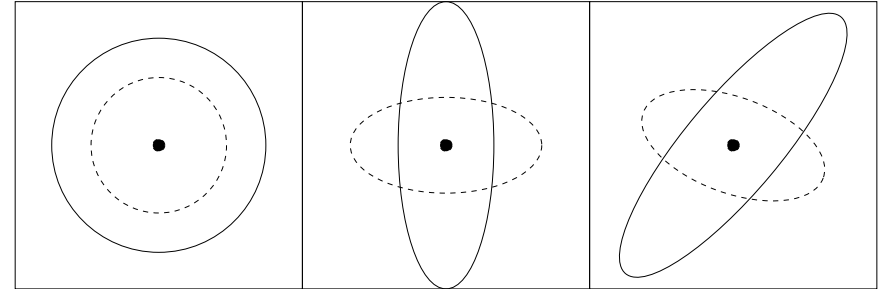
33

---

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

Lines of Equal Density



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

68%, 95%, 99.7% of samples drop into $(x - m)^{\mathrm{T}} \mathbf{C}^{-1}(x - m) \lessgtr n - \frac{1}{2} \pm \frac{1^2}{2}$, $\pm \frac{2^2}{2}, \pm \frac{3^2}{2}$

... ESs

34

---

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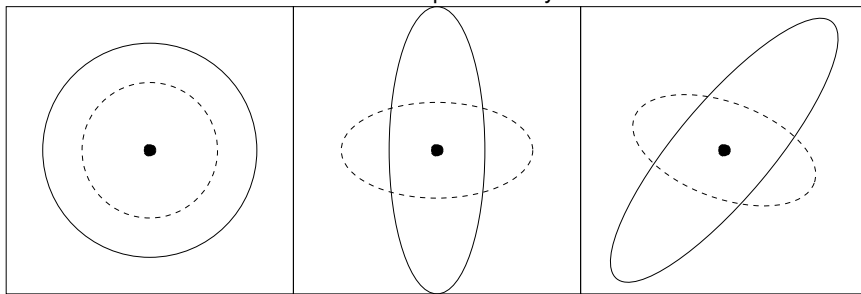
Lines of Equal Density



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

68%, 95%, 99.7% of samples drop into $(x - m)^{\mathrm{T}} \mathbf{C}^{-1}(x - m) \lessgtr n - \frac{1}{2} \pm \frac{1^2}{2}$, $\pm \frac{2^2}{2}, \pm \frac{3^2}{2}$

... ESs

34

---

# Evolution Strategies

Terminology

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

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

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

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

Sample one offspring from parent $m$

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

If $x$ better than $m$ select

$$m \leftarrow x$$

... why?

35

---

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

Non-elitist selection and intermediate (weighted) recombination

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

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

where

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

36

---

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

Non-elitist selection and intermediate (weighted) recombination

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

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.

36

---

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

Non-elitist selection and intermediate (weighted) recombination

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

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

$$w_1 \geq \cdots \geq w_\mu > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4}$$
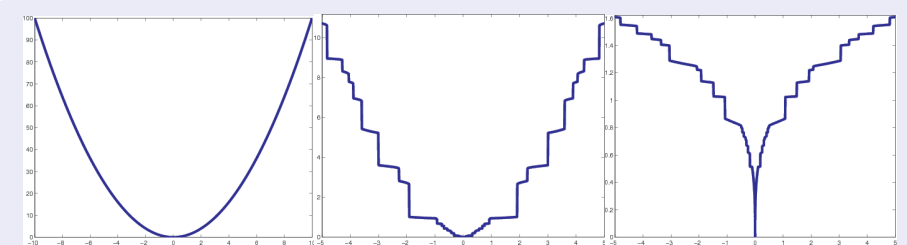
The best $\mu$ points are selected from the new solutions (non-elitistic)
and weighted intermediate recombination is applied.

36

---

## Invariance Under Monotonically Increasing Functions

### Rank-based algorithms

Update of all parameters uses only the ranks

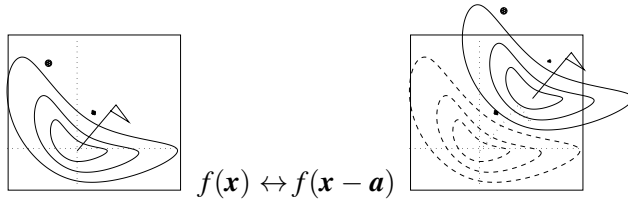$$f(x_{1:\lambda}) \leq f(x_{2:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda})$$



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

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

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

37

# 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 $f_a$

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

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

38

---

# Summary

On 20D Sphere Function: $f(\mathbf{x}) = \sum_{i=1}^{N} [\mathbf{x}]_i^2$
- ES without adaptation can't approach the optimum $\Rightarrow$ adaptation required
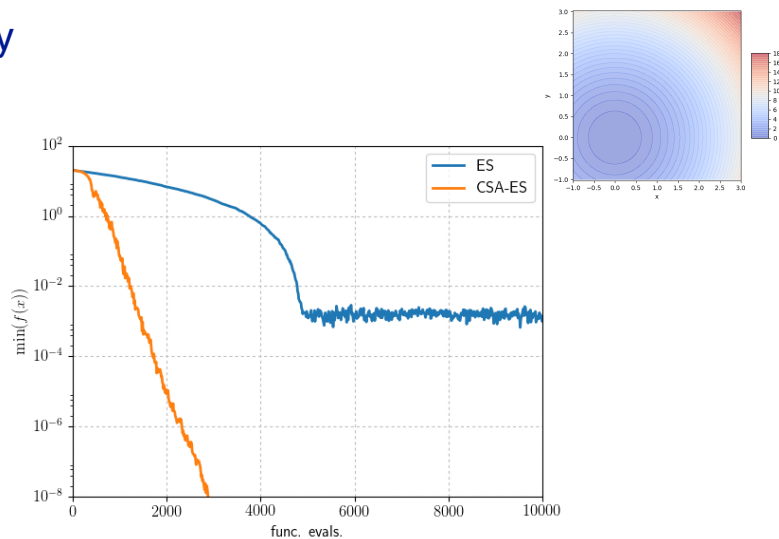
39

---

# Summary

On 20D Sphere Function: $f(\mathbf{x}) = \sum_{i=1}^{N} [\mathbf{x}]_i^2$
- ES without adaptation can't approach the optimum $\Rightarrow$ adaptation required

39

---

# Evolution Strategies

Recalling

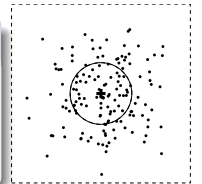### New search points are sampled normally distributed

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

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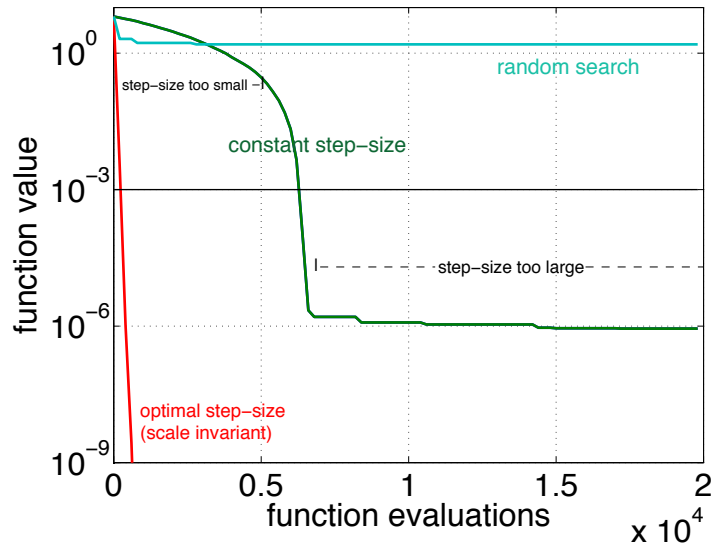
where

- the mean vector $\boldsymbol{m} \in \mathbb{R}^n$ represents the favorite solution and $\boldsymbol{m} \leftarrow \sum_{i=1}^{\mu} w_i \, \boldsymbol{x}_{i:\lambda}$
- 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

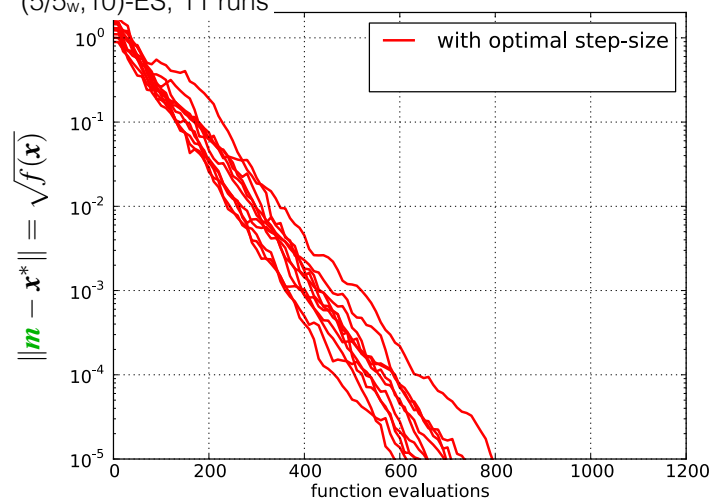The remaining question is how to update $\sigma$ and $\mathbf{C}$.

40

# 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 $n = 10$

41

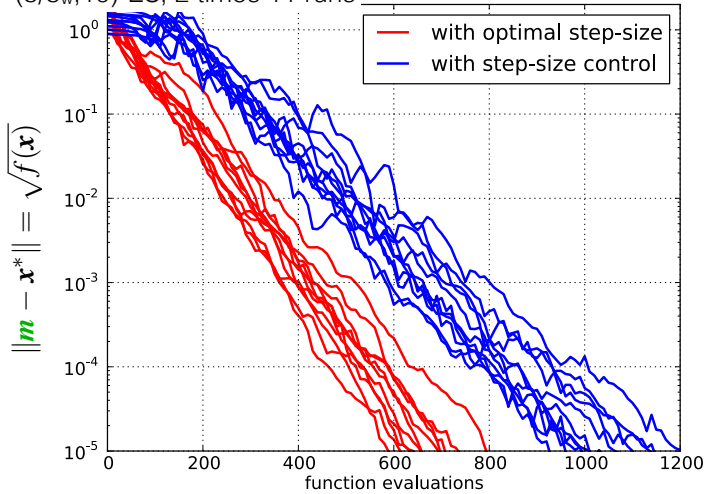# Why Step-Size Control?

(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 $\sigma$

42

# 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 $\sigma$

43

# 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 $\|\boldsymbol{m}\| = 10^{-5}$: $\frac{1100-100}{650} \approx 1.5$

44

## Why Step-Size Control?
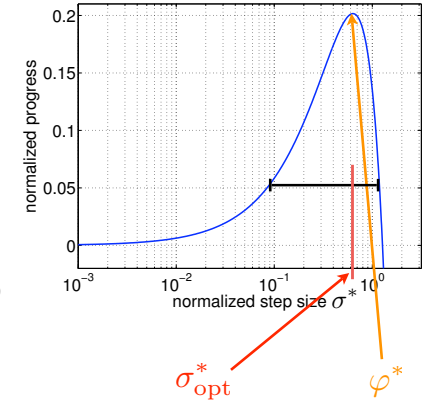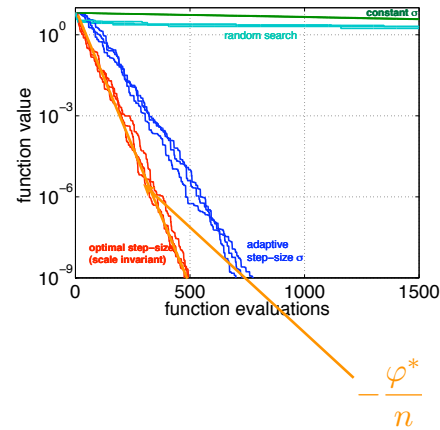
$(5/5_w,10)$-ES



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

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

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

45

## 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

46

## 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*
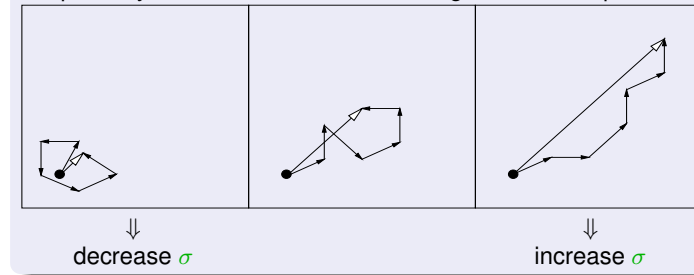
47

## 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)

48

## Path Length Control (CSA)
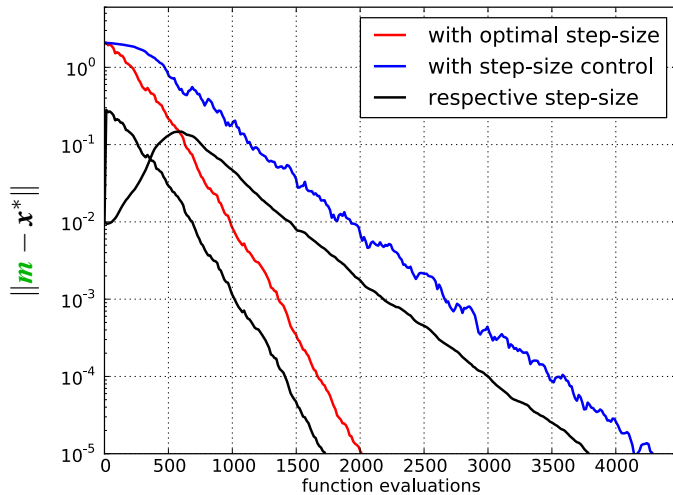The Equations

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

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

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

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

49

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



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

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

50

## Step-Size Control: Summary



Why Step-Size Control?
- to achieve linear convergence

Cumulative Step-Size Adaptation
- efficient and robust for $\lambda \leq N$
- inefficient (1) $\lambda \gg N$, (2) function with ineffective axes

Alternative Step-Size Adaptation Mechanisms
- Two-Point Step-Size Adaptation
- Median Success Rule, Population Success Rule

*the effective adaptation of the overall population diversity seems yet to pose open questions, in particular with recombination or without entire control over the realized distribution.[a]*

[a]Hansen et al. How to Assess Step-Size Adaptation Mechanisms in Randomised Search. PPSN 2014

51

# Step-Size Control: Summary



On 20D TwoAxes Function: $f(\mathbf{x}) = \sum_{i=1}^{N/2}[\mathbf{Rx}]_i^2 + a^2 \sum_{i=N/2+1}^{N}[\mathbf{Rx}]_i^2$, $\mathbf{R}$: orthogonal
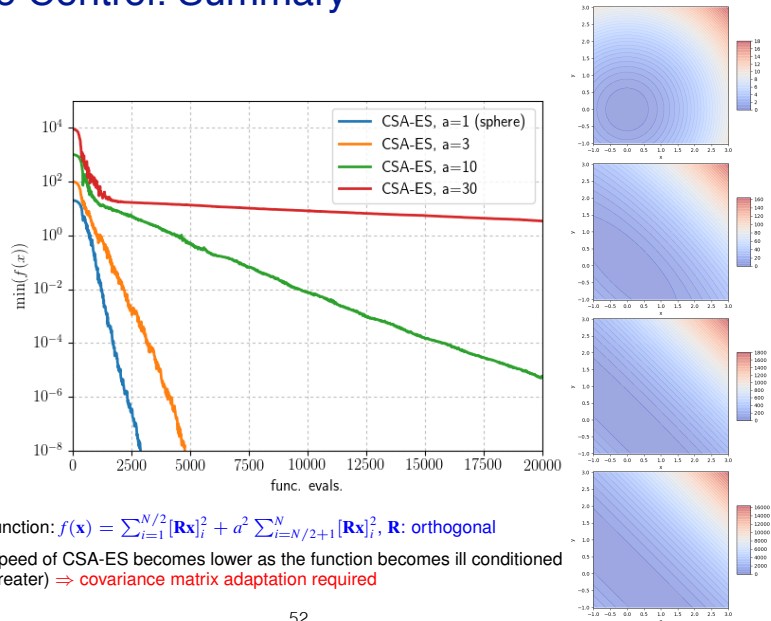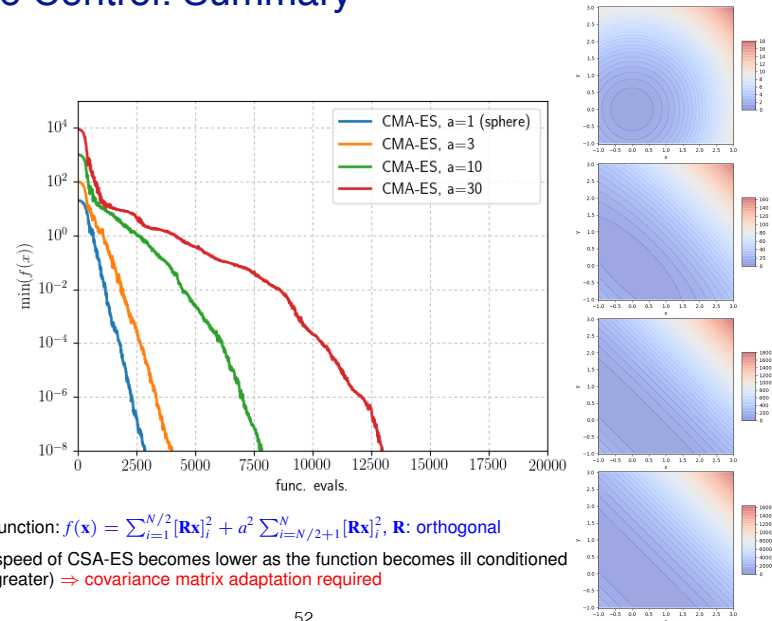
- convergence speed of CSA-ES becomes lower as the function becomes ill conditioned ($a^2$ becomes greater) ⇒ covariance matrix adaptation required

52

# Evolution Strategies
Recalling

**New search points are sampled normally distributed**

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

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 $\mathbf{C}$.

53

# Covariance Matrix Adaptation
Rank-One Update

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

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

...equations

54

658

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

55

# 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 $\sigma$)

56

# 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 $y_w$,
$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times y_w y_w^{\mathrm{T}}$

57

# 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 $\sigma$)

58

# 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 $\sigma$)
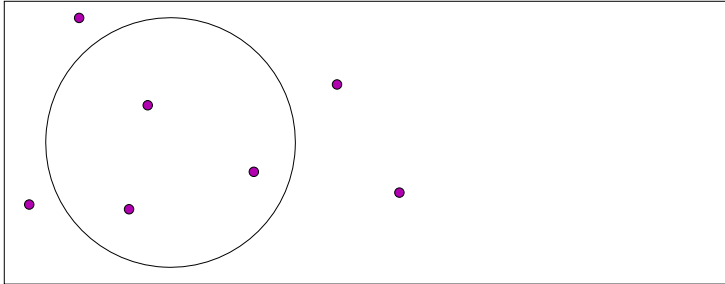
---

# 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 $y_w$,
$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times y_w y_w^{\mathrm{T}}$$

---

# 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

# Covariance Matrix Adaptation
## Rank-One Update

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

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

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

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

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
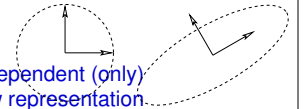[8] Ljung 1999. System Identification: Theory for the User
[9] Haario et al 2001. An adaptive Metropolis algorithm, JSTOR

63

---

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

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

cumulation—rank

64

---

# 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

$$p_{\mathbf{c}} \propto \sum_{i=0}^{g} \underbrace{(1 - c_{\mathbf{c}})^{g-i}}_{\substack{\text{exponentially} \\ \text{fading weights}}} y_w^{(i)}$$

The recursive construction of the evolution path (cumulation):

$$p_{\mathbf{c}} \leftarrow \underbrace{(1 - c_{\mathbf{c}})\, p_{\mathbf{c}}}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_{\mathbf{c}})^2}\sqrt{\mu_w}}_{\text{normalization factor}}\; \underbrace{y_w}_{\text{input} = \frac{m - m_{\text{old}}}{\sigma}}$$

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

65

---

# 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

$$p_{\mathbf{c}} \propto \sum_{i=0}^{g} \underbrace{(1 - c_{\mathbf{c}})^{g-i}}_{\substack{\text{exponentially} \\ \text{fading weights}}} y_w^{(i)}$$

The recursive construction of the evolution path (cumulation):

$$p_{\mathbf{c}} \leftarrow \underbrace{(1 - c_{\mathbf{c}})\, p_{\mathbf{c}}}_{\text{decay factor}} + \underbrace{\sqrt{1 - (1 - c_{\mathbf{c}})^2}\sqrt{\mu_w}}_{\text{normalization factor}}\; \underbrace{y_w}_{\text{input} = \frac{m - m_{\text{old}}}{\sigma}}$$

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

66

"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?

67

---

# Cumulation

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

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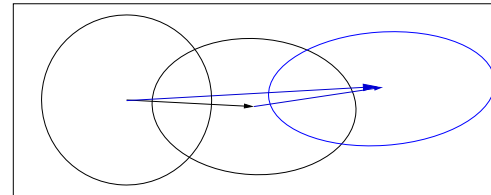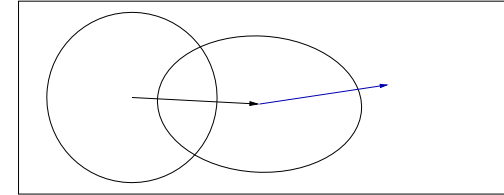We used $\mathbf{y}_w \mathbf{y}_w^{\text{T}}$ for updating $\mathbf{C}$. Because $\mathbf{y}_w \mathbf{y}_w^{\text{T}} = -\mathbf{y}_w(-\mathbf{y}_w)^{\text{T}}$ the sign of $\mathbf{y}_w$ is lost.



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

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

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

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".

68

---

# Cumulation

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

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



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

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

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

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".

69

---

# Cumulation

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

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



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

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

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

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".
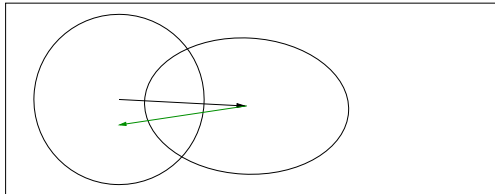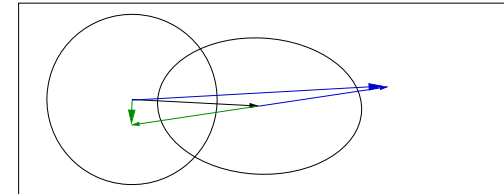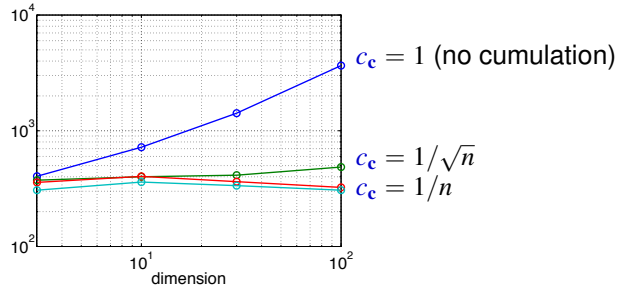
70

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 $f(\boldsymbol{x}) = x_1^2 + 10^6 \sum_{i=2}^{n} x_i^2$



$c_{\mathbf{c}} = 1$ (no cumulation)

$c_{\mathbf{c}} = 1/\sqrt{n}$
$c_{\mathbf{c}} = 1/n$

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

71

---

# Rank-$\mu$ Update

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

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 $c_{\mathrm{cov}} \le 1$.

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

72

---

# Rank-$\mu$ Update

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

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 $c_{\mathrm{cov}} \le 1$.

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

72

---

# Rank-$\mu$ Update

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

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]

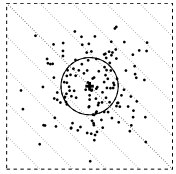The rank-$\mu$ update then reads

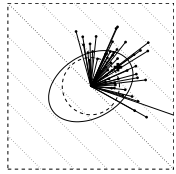$$\mathbf{C} \leftarrow (1 - c_{\mathrm{cov}})\,\mathbf{C} + c_{\mathrm{cov}}\,\mathbf{C}_\mu$$

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

[10]Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.
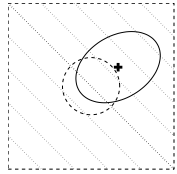
72

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

$$C_\mu = \frac{1}{\mu} \sum y_{i:\lambda} y_{i:\lambda}^T$$
$$C \leftarrow (1-1) \times C + 1 \times C_\mu$$

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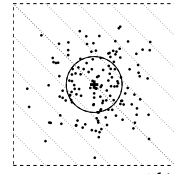
new distribution

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

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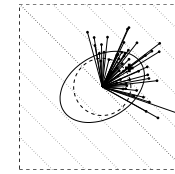$$w_1 = \cdots = w_\mu = \frac{1}{\mu},$$
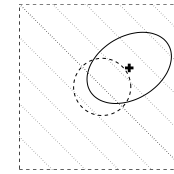and $c_{cov} = 1$

---

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

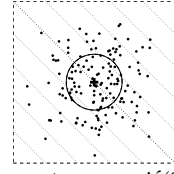

$$x_i = m_{old} + y_i, \quad y_i \sim \mathcal{N}(0, C)$$

$$C \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - m_{old})(x_{i:\lambda} - m_{old})^T$$

$$m_{new} = m_{old} + \frac{1}{\mu} \sum y_{i:\lambda}$$

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

$$x_i = m_{old} + y_i, \quad y_i \sim \mathcal{N}(0, C)$$

$$C \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - m_{new})(x_{i:\lambda} - m_{new})^T$$

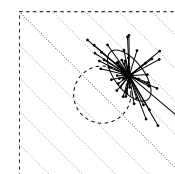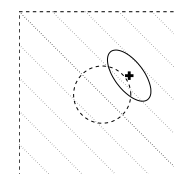$$m_{new} = m_{old} + \frac{1}{\mu} \sum y_{i:\lambda}$$

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

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

calculating $C$ from $\mu = 50$ solutions

new distribution

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

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 $\lambda \geq 3n + 10$

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 $\lambda \geq 3n + 10$

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 $\mathcal{O}(n)$ .
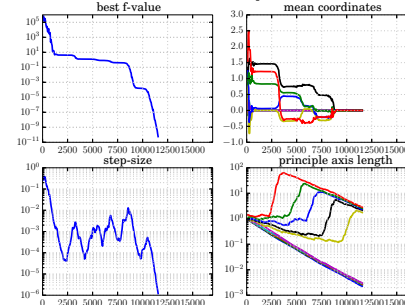
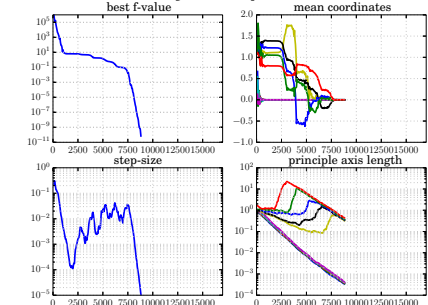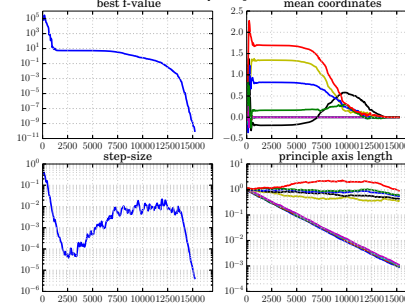Rank-one update and rank-$\mu$ update can be combined

...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
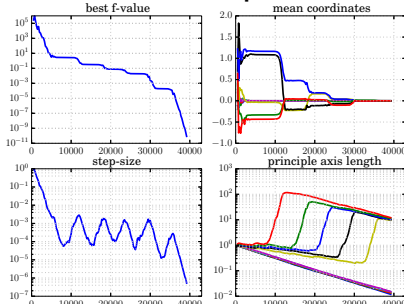
---

**Rank-one update**

**Hybrid update**

**Rank-$\mu$ update**

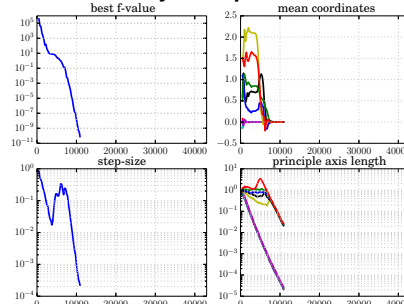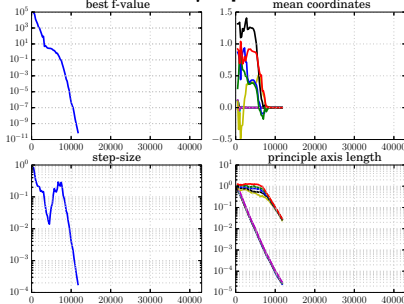$$f_{\text{TwoAxes}}(x) = \sum_{i=1}^{5} x_i^2 + 10^6 \sum_{i=6}^{10} x_i^2$$

$\lambda = 10$ (default for $N = 10$)

---

**Rank-one update**

**Hybrid update**

**Rank-$\mu$ update**

$$f_{\text{TwoAxes}}(x) = \sum_{i=1}^{5} x_i^2 + 10^6 \sum_{i=6}^{10} x_i^2$$

$\lambda = 50$

---

# Summary of Equations

The Covariance Matrix Adaptation Evolution Strategy

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

Initialize: $\mathbf{C} = \mathbf{I}$, and $\boldsymbol{p}_\mathbf{c} = \mathbf{0}$, $\boldsymbol{p}_\sigma = \mathbf{0}$,

Set: $c_\mathbf{c} \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n^2$, $c_\mu \approx \mu_w/n^2$, $c_1 + c_\mu \leq 1$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$,

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

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

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

$$\boldsymbol{p}_\mathbf{c} \leftarrow (1 - c_\mathbf{c})\,\boldsymbol{p}_\mathbf{c} + \mathbb{1}_{\{\|p_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_\mathbf{c})^2}\sqrt{\mu_w}\,\boldsymbol{y}_w \qquad \text{cumulation for } \mathbf{C}$$

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

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

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

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

## 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 $w_{i=1...\lambda}$ 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

80

---

## Topics

1. What makes the problem difficult to solve?

2. How does the CMA-ES work?

- Normal Distribution, Rank-Based Recombination
- Step-Size Adaptation (CSA)
- Covariance Matrix Adaptation (Hybrid-CMA)

### 3. What can/should the users do for the CMA-ES to work effectively on your problem?

- Restart, Increasing Population Size
- Restricted Covariance Matrix

81

---

## Default Parameter Values

CMA-ES + (B)IPOP Restart Strategy = Quasi-Parameter Free Optimizer

The following parameters were identified in carefully chosen experimental set ups.

- related to selection and recombination
  - $\lambda$: offspring number, new solutions sampled, population size
  - $\mu$: parent number, solutions involved in updates of
  - $w_i$: recombination weights
- related to $C$-update
  - $c_c$: decay rate for the evolution path, cumulation factor
  - $c_1$: learning rate for rank-one update of $C$
  - $c_\mu$: learning rate for rank-$\mu$ update of $C$
- related to $\sigma$-update
  - $c_\sigma$: decay rate of the evolution path
  - $d_\sigma$: damping for $\sigma$-change

82

---

## Default Parameter Values

CMA-ES + (B)IPOP Restart Strategy = Quasi-Parameter Free Optimizer

The following parameters were identified in carefully chosen experimental set ups.

- related to selection and recombination
  - $\lambda$: offspring number, new solutions sampled, population size
  - $\mu$: parent number, solutions involved in updates of
  - $w_i$: recombination weights
- related to $C$-update
  - $c_c$: decay rate for the evolution path, cumulation factor
  - $c_1$: learning rate for rank-one update of $C$
  - $c_\mu$: learning rate for rank-$\mu$ update of $C$
- related to $\sigma$-update
  - $c_\sigma$: decay rate of the evolution path
  - $d_\sigma$: damping for $\sigma$-change

The default values depends only on dimension N. They do in the first place not depend on the objective function.

82

# Parameters to be set depending on the problem
Initialization and termination conditions

The following should be set or implemented depending on the problem.

- related to the initial search distribution
  - $m^{(0)}$: initial mean vector
  - $\sigma^{(0)}$ (or $\sqrt{C_{i,i}^{(0)}}$): initial (coordinate-wise) standard deviation
- related to stopping conditions
  - max. func. evals.
  - max. iterations
  - function value tolerance
  - min. axis length
  - stagnation

83

---

# Parameters to be set depending on the problem
Initialization and termination conditions

The following should be set or implemented depending on the problem.

- related to the initial search distribution
  - $m^{(0)}$: initial mean vector
  - $\sigma^{(0)}$ (or $\sqrt{C_{i,i}^{(0)}}$): initial (coordinate-wise) standard deviation
- related to stopping conditions
  - max. func. evals.
  - max. iterations
  - function value tolerance
  - min. axis length
  - stagnation

Given an initial search interval $[a_i, b_i]$ for $i = 1, \ldots, n$, a reasonable choice will be

- $m_i^{(0)} = (a_i + b_i)/2$   or   $m_i^{(0)} \sim \mathcal{U}[a_i + \epsilon, b_i - \epsilon]$
- $\sqrt{C_{i,i}^{(0)}} = \frac{b_i - a_i}{2 \text{ to } 4}$ for $i = 1, \ldots, n$ and $C_{i,j}^{(0)} = 0$ for $i \neq j$

83

---

# Python CMA-ES Implementation
https://github.com/CMA-ES/pycma

## pycma

A Python implementation of CMA-ES and a few related numerical optimization tools.

The Covariance Matrix Adaptation Evolution Strategy (CMA-ES) is a stochastic numerical optimization algorithm for difficult (non-convex, ill-conditioned, multi-modal, rugged, noisy) optimization problems in continuous search spaces.

The API Documentation is available here.

## Installation

Download and unzip the code (see green button above) or `git clone https://github.com/CMA-ES/pycma.git` .

- Either, copy (or move) the `cma` source code folder into a folder visible to Python, namely a folder which is in the Python path (e.g. the current folder). Then, `import cma` works without any further installation.
- Or, install the `cma` package by typing within the folder, where the `cma` source code folder is visible,

  ```
  python -m pip install -e cma
  ```

Typing `pip` instead of `python -m pip` may be sufficient, prefixing with `sudo` may be necessary. Moving the `cma` folder away from this location would invalidate the installation.

84

---

# Python CMA-ES Demo
https://github.com/CMA-ES/pycma

Optimizing 15D Tablet Function

```
import cma
opts = cma.CMAOptions()
opts['tolfun'] = 1e-4          # f-tolerance
opts['ftarget'] = 1e-4         # f-target value
opts['maxfevals'] = 1e6        # max #FEs
# opts['popsize'] = '10 * N' # population size
es = cma.CMAEvolutionStrategy(x0=15 * [1],   # Initial mean vector
                              sigma0=1,       # Initial step-size
                              inopts=opts     # Options
                              ).optimize(cma.ff.tablet) # Objective
```

```
(6_w,12)-aCMA-ES (mu_w=3.7,w_1=40%) in dimension 15 (seed=137090, Mon Apr 24 14:58:52 2017)
Iterat #Fevals   function value  axis ratio  sigma  min&max std  t[m:s]
    1      12 1.537676704740862e+02 1.0e+00 1.03e+00  1e+00  1e+00 0:00.0
    2      24 1.408854302050177e+02 1.1e+00 1.03e+00  1e+00  1e+00 0:00.0
    3      36 3.712560411998829e+03 1.2e+00 1.02e+00  1e+00  1e+00 0:00.0
  100    1200 1.506902133117476e+02 1.7e+01 5.06e-01  6e-02  7e-01 0:00.1
  200    2400 1.893840652870748e+01 1.9e+02 2.99e-01  3e-03  4e-01 0:00.3
  300    3600 3.434648669054741e-01 7.1e+02 1.30e-01  3e-04  1e-01 0:00.4
  384    4608 9.403602672438789e-05 1.1e+03 4.69e-03  4e-06  3e-03 0:00.5
```

From a practical perspective: given an unknown optimisation problem, the first thing I tend to do is try to improve a given (initial) solution using a small initial sigma. Then I (can) increase sigma successively (by a factor of 10 or more, depending on what I have seen in the initial evolution of sigma previously) and see whether I find the same or better (or worse) solutions.
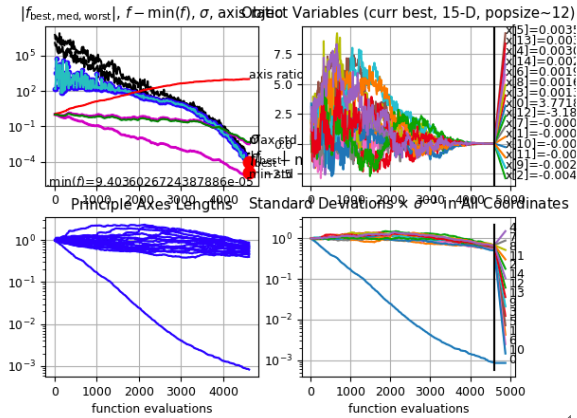
85

# Python CMA-ES Demo

https://github.com/CMA-ES/pycma

Optimizing 15D Tablet Function



---

# Multimodality

---

# Multimodality

Two approaches for multimodal functions: Try again with

---

# Multimodality

Two approaches for multimodal functions: Try again with
- a larger population size

# Multimodality

Two approaches for multimodal functions: Try again with
- a larger population size
- a smaller initial step-size (and random initial mean vector)

87

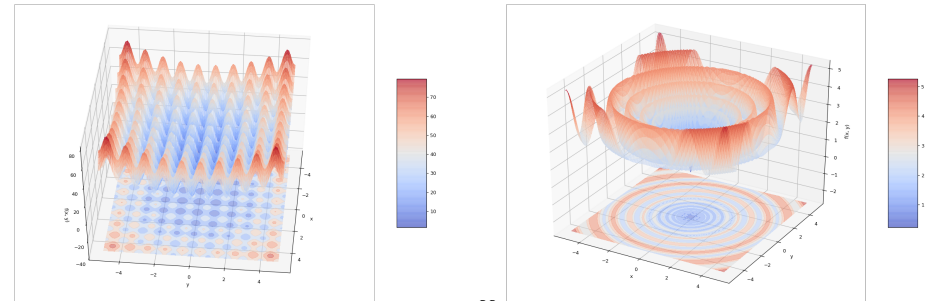# Multimodality

Two approaches for multimodal functions: Try again with
- a larger population size
- a smaller initial step-size (and random initial mean vector)

A restart with a large population size helps if the objective function has a well global structure
- functions such as Schaffer, Rastrigin, BBOB function 15~19
- loosely, unimodal global structure + deterministic noise

88

# Multimodality

Hansen and Kern. Evaluating the CMA Evolution Strategy on Multimodal Test Functions, PPSN 2004.
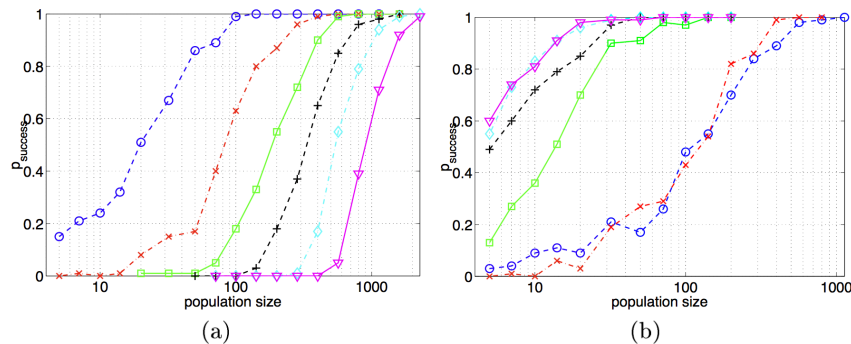
**Fig. 1.** Success rate to reach $f_{\text{stop}} = 10^{-10}$ versus population size for (a) Rastrigin function (b) Griewank function for dimensions $n = 2$ ('$--\bigcirc--$'), $n = 5$ ('$-\cdot-\times-\cdot-$'), $n = 10$ ('$-\square-$'), $n = 20$ ('$--+--$'), $n = 40$ ('$-\cdot-\diamondsuit-\cdot-$'), and $n = 80$ ('$-\bigtriangledown-$').
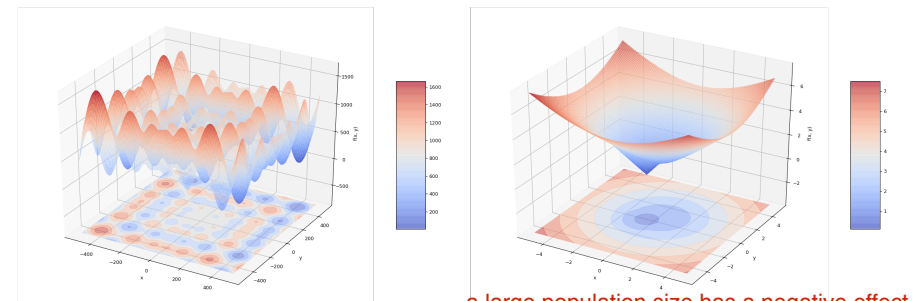
89

# Multimodality

Two approaches for multimodal functions: Try again with
- a larger population size
- a smaller initial step-size (and random initial mean vector)

A restart with a small initial step-size helps if the objective function has a weak global structure
- functions such as Schwefel, Bi-Sphere, BBOB function 20~24

a large population size has a negative effect

90

# Restart Strategy
It makes the CMA-ES parameter free

IPOP: Restart with increasing the population size
- start with the default population size
- double the population size after each trial (parameter sweep)
- may be considered as gold standard for automated restarts

BIPOP: IPOP regime + Local search regime
- IPOP regime: restart with increasing population size
- Local search regime: restart with a smaller step-size and a smaller population size than the IPOP regime

91

---

# Topics

3. What can/should the users do for the CMA-ES to work efficiently on your problem?

- Restart, Increasing Population Size
- Restricted Covariance Matrix

92

---

# Motivation of the Restricted Covariance Matrix

Bottlenecks of the CMA-ES on high dimensional problems
1. $\mathcal{O}(N^2)$ Time and Space Complexities
   - to store and update $C \in \mathbb{R}^{N \times N}$
   - to compute the eigen decomposition of $C$
2. $\mathcal{O}(1/N^2)$ Learning Rates for $C$-Update
   - $c_\mu \approx \mu_w/N^2$
   - $c_1 \approx 2/N^2$

Exploit prior knowledge on the problem structure such as separability

$\Rightarrow$ decrease the degrees of freedom of the covariance matrix for
- less time and space complexities
- a higher learning rates that potentially accelerate the adaptation

93

---

# Motivation of the Restricted Covariance Matrix

Bottlenecks of the CMA-ES on high dimensional problems
1. $\mathcal{O}(N^2)$ Time and Space Complexities
   - to store and update $C \in \mathbb{R}^{N \times N}$
   - to compute the eigen decomposition of $C$
2. $\mathcal{O}(1/N^2)$ Learning Rates for $C$-Update
   - $c_\mu \approx \mu_w/N^2$
   - $c_1 \approx 2/N^2$

Exploit prior knowledge on the problem structure such as separability

$\Rightarrow$ decrease the degrees of freedom of the covariance matrix for
- less time and space complexities
- a higher learning rates that potentially accelerate the adaptation

93

# Motivation of the Restricted Covariance Matrix

Bottlenecks of the CMA-ES on high dimensional problems

1. $\mathcal{O}(N^2)$ Time and Space Complexities
   - to store and update $C \in \mathbb{R}^{N \times N}$
   - to compute the eigen decomposition of $C$
2. $\mathcal{O}(1/N^2)$ Learning Rates for $C$-Update
   - $c_\mu \approx \mu_w / N^2$
   - $c_1 \approx 2/N^2$

Exploit prior knowledge on the problem structure such as separability

$\Rightarrow$ decrease the degrees of freedom of the covariance matrix for
- less time and space complexities
- a higher learning rates that potentially accelerate the adaptation
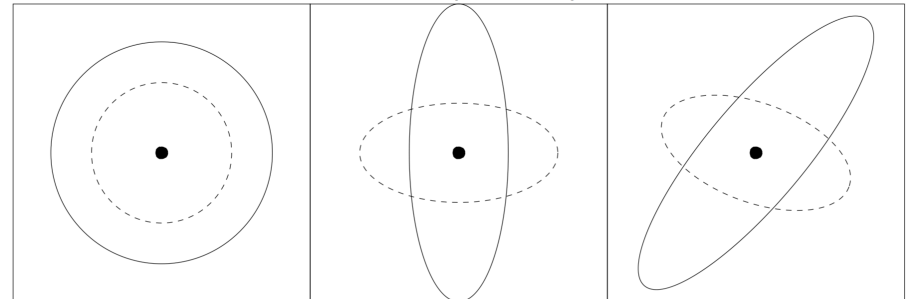
93

---

---

# Variants with Restricted Covariance Matrix

CMA-ES Variants with Restricted Covariance Matrices
- Sep-CMA [Ros and Hansen, 2008]
  - $C = \mathbf{D}$. $\mathbf{D}$: Diagonal
- VD-CMA [Akimoto et al., 2014]
  - $C = \mathbf{D}(\mathbf{I} + \boldsymbol{v}\boldsymbol{v}^{\mathrm{T}})\mathbf{D}$. $\mathbf{D}$: Diagonal, $\boldsymbol{v} \in \mathbb{R}^N$.
- LM-CMA [Loshchilov, 2014]
  - $C = \mathbf{I} + \sum_{i=1}^k \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{T}}$. $\boldsymbol{v}_i \in \mathbb{R}^N$.
- VkD-CMA [Akimoto and Hansen, 2016]
  - $C = \mathbf{D}(\mathbf{I} + \sum_{i=1}^k \boldsymbol{v}_i\boldsymbol{v}_i^{\mathrm{T}})\mathbf{D}$. $\boldsymbol{v}_i \in \mathbb{R}^N$.

[Ros and Hansen, 2008] Ros, R. and Hansen, N. (2008). A simple modification in CMA-ES achieving linear time and space complexity. In Parallel Problem Solving from Nature - PPSN X, pages 296–305. Springer.

[Akimoto et al., 2014] Akimoto, Y., Auger, A., and Hansen, N. (2014). Comparison-based natural gradient optimization in high dimension. In Proceedings of Genetic and Evolutionary Computation Conference, pages 373–380, Vancouver, BC, Canada.

[Loshchilov, 2014] Loshchilov, I. (2014). A computationally efficient limited memory cma-es for large scale optimization. In Proceedings of Genetic and Evolutionary Computation Conference, pages 397–404.

[Akimoto and Hansen, 2016] Akimoto, Y. and Hansen, N. (2016). Projection-based restricted covariance matrix adaptation for high dimension. In Genetic and Evolutionary Computation Conference, GECCO 2016, Denver, Colorado, USA, July 20-24, 2016, page (accepted). ACM.
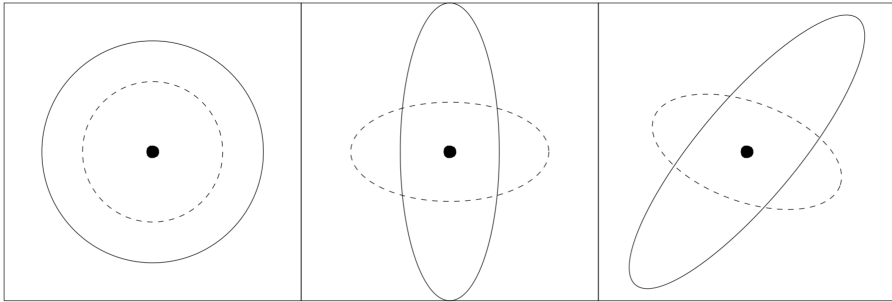
94

---

# Separable CMA (Sep-CMA)



$\mathcal{N}(\boldsymbol{m}, \sigma^2\mathbf{I}) \sim \boldsymbol{m} + \sigma\mathcal{N}(\mathbf{0}, \mathbf{I})$
one degree of freedom $\sigma$

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

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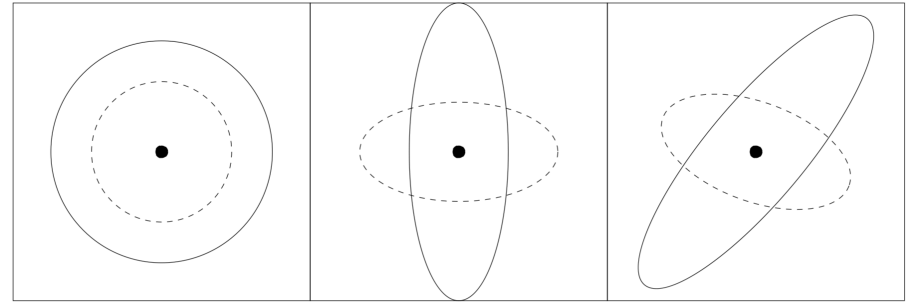$(n^2 + n)/2$ degrees of freedom

95

## Separable CMA (Sep-CMA)

$\mathcal{N}\left(m, \sigma^2 \mathbf{I}\right) \sim m + \sigma \mathcal{N}(0, \mathbf{I})$
one degree of freedom $\sigma$

$\mathcal{N}\left(m, \mathbf{D}^2\right) \sim m + \mathbf{D}\,\mathcal{N}(0, \mathbf{I})$
$n$ degrees of freedom

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

CMA $\quad C_{\text{cma}}^{(t+1)} = C^{(t)} + c_1\left(p_c\,p_c^{\text{T}} - C^{(t)}\right) + c_\mu \sum_{i=1}^{\mu} w_i\left((x_i - m^{(t)})(x_i - m^{(t)})^{\text{T}} - C^{(t)}\right)$

SEP $\quad [C_{\text{sep}}^{(t+1)}]_{k,k} = [C^{(t)}]_{k,k} + c_1\left([p_c]_k^2 - [C^{(t)}]_{k,k}\right) + c_\mu \sum_{i=1}^{\mu} w_i\left([x_i - m^{(t)}]_k^2 - [C^{(t)}]_{k,k}\right)$

95

---

## Separable CMA (Sep-CMA)

$\mathcal{N}\left(m, \sigma^2 \mathbf{I}\right) \sim m + \sigma \mathcal{N}(0, \mathbf{I})$
one degree of freedom $\sigma$

$\mathcal{N}\left(m, \mathbf{D}^2\right) \sim m + \mathbf{D}\,\mathcal{N}(0, \mathbf{I})$
$n$ degrees of freedom

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

CMA $\quad C_{\text{cma}}^{(t+1)} = C^{(t)} + c_1\left(p_c\,p_c^{\text{T}} - C^{(t)}\right) + c_\mu \sum_{i=1}^{\mu} w_i\left((x_i - m^{(t)})(x_i - m^{(t)})^{\text{T}} - C^{(t)}\right)$

SEP $\quad [C_{\text{sep}}^{(t+1)}]_{k,k} = [C^{(t)}]_{k,k} + c_1\left([p_c]_k^2 - [C^{(t)}]_{k,k}\right) + c_\mu \sum_{i=1}^{\mu} w_i\left([x_i - m^{(t)}]_k^2 - [C^{(t)}]_{k,k}\right)$

→ (N + 2)/3 times greater than CMA
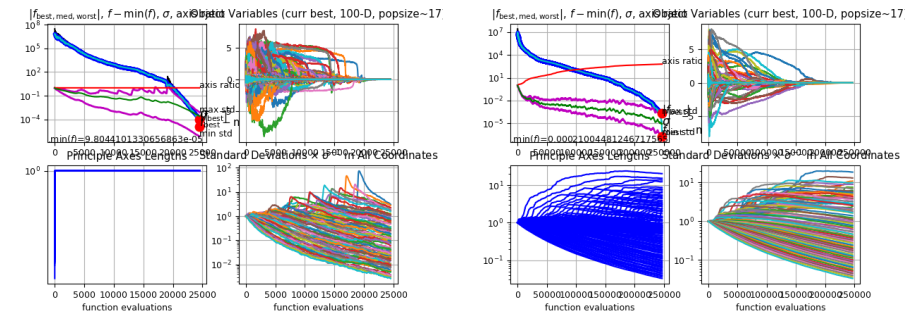
95

---

## Demo: On 100D Separable Ellipsoid Function

Separable-CMA

CMA

96

---

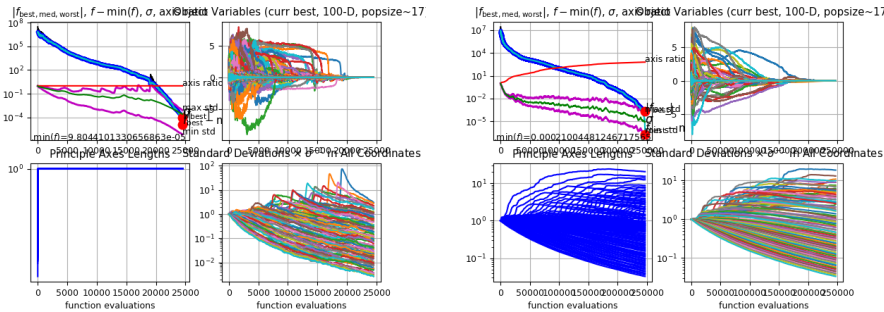## Demo: On 100D Separable Ellipsoid Function

Separable-CMA

CMA

- CMA needed 10 times more FEs + more CPU time

96

# Demo: On 100D Separable Ellipsoid Function



Separable-CMA



CMA

- CMA needed 10 times more FEs + more CPU time
- However, Sep-CMA won't be able to solve rotated ellipsoid function as efficiently as it solves separable ellipsoid

96

---

# Summary and Final Remarks

97

---

# 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

98

---

# Main Characteristics of (CMA) Evolution Strategies

1. Multivariate normal distribution to generate new search points

   follows the maximum entropy principle

2. Rank-based selection

   implies invariance, same performance on $g(f(x))$ for any increasing $g$
   more invariance properties are featured

3. Step-size control facilitates fast (log-linear) convergence and possibly linear scaling with the dimension

   in CMA-ES based on an evolution path (a non-local trajectory)

4. *Covariance matrix adaptation (CMA)* increases the likelihood of previously successful steps and can improve performance by orders of magnitude

   the update follows the natural gradient
   $\mathbf{C} \propto \boldsymbol{H}^{-1} \iff$ adapts a variable metric
   $\iff$ new (rotated) problem representation
   $\implies f : \boldsymbol{x} \mapsto g(\boldsymbol{x}^\mathsf{T} \boldsymbol{H} \boldsymbol{x})$ reduces to $\boldsymbol{x} \mapsto \boldsymbol{x}^\mathsf{T} \boldsymbol{x}$

99

## 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

  variants with restricted covariance matrix such as Sep-CMA

- 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

100

# Thank you

Source code for CMA-ES in C, C++, Java, Matlab, Octave, Python, R, Scilab is available (or linked to) at
http://cma.gforge.inria.fr/cmaes_sourcecode_page.html

101