

Tutorial—Evolution Strategies and Related Estimation of Distribution Algorithms

Anne Auger & Nikolaus Hansen

INRIA Saclay - Ile-de-France, project team TAO
Universite Paris-Sud, LRI, Bat. 490
91405 ORSAY Cedex, France

Copyright is held by the author/owner(s).
GECCO'08, July 12, 2008, Atlanta, Georgia, USA.
ACM 978-1-60558-131-6/08/07

Content

- 1 Problem Statement
 - Black Box Optimization and Its Difficulties
 - Non-Separable Problems
 - Ill-Conditioned Problems
- 2 Evolution Strategies and EDAs
 - A Search Template
 - The Normal Distribution
- 3 Step-Size Control
 - Why Step-Size Control
 - One-Fifth Success Rule
 - Self-Adaptation
 - Path Length Control
- 4 Covariance Matrix Adaptation
 - Covariance Matrix Rank-One Update
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
 - Estimation of Distribution
- 5 Conclusion

Problem Statement

Continuous Domain Search/Optimization

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

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

- **Black Box** scenario (direct search scenario)



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

Problem Statement

Continuous Domain Search/Optimization

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

Approach: stochastic search, Evolutionary Algorithms

Metaphors

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

... function properties

Objective Function Properties

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

Additionally, f can be

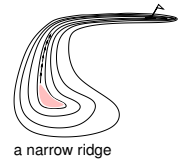
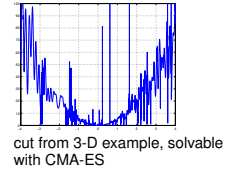
- multimodal there are eventually many local optima
- non-smooth derivatives do not exist
- discontinuous
- ill-conditioned
- noisy
- ...

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

What Makes a Function Difficult to Solve?

Why stochastic search?

- ruggedness
non-smooth, discontinuous, multimodal, and/or noisy function
- dimensionality
(considerably) larger than three
- non-separability
dependencies between the objective variables
- ill-conditioning



Separable Problems

Definition (Separable Problem)

A function f is separable if

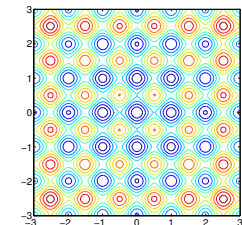
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left(\arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

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

Example: Additively decomposable functions

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

Rastrigin function



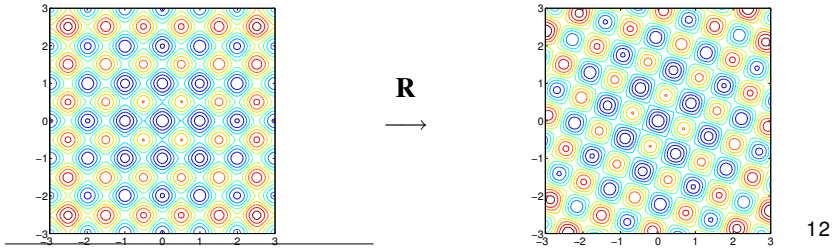
Non-Separable Problems

Building a non-separable problem from a separable one

Rotating the coordinate system

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

R rotation matrix



12

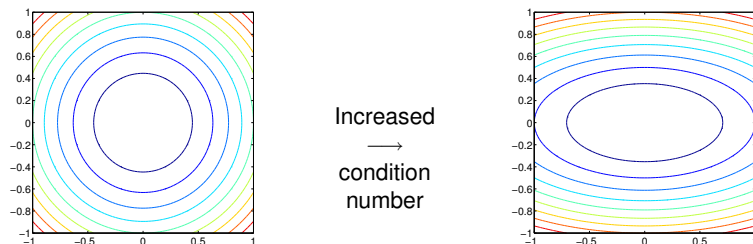
¹Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

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

Ill-Conditioned Problems

If f is quadratic, $f : x \mapsto x^T H x$, ill-conditioned means a high condition number of Hessian Matrix H

ill-conditioned means "squeezed" lines of equal function value



consider the curvature of iso-fitness lines

What Makes a Function Difficult to Solve?

... and what can be done

The Problem	What can be done
Ruggedness	non-local policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed
Dimensionality, Non-Separability	stochastic, non-elitistic, population-based method recombination operator serves as repair mechanism
Ill-conditioning	exploiting the problem structure locality, neighborhood, encoding
	second order approach changes the neighborhood metric

- 1 Problem Statement
 - Black Box Optimization and Its Difficulties
 - Non-Separable Problems
 - Ill-Conditioned Problems
- 2 Evolution Strategies and EDAs
 - A Search Template
 - The Normal Distribution
- 3 Step-Size Control
 - Why Step-Size Control
 - One-Fifth Success Rule
 - Self-Adaptation
 - Path Length Control
- 4 Covariance Matrix Adaptation
 - Covariance Matrix Rank-One Update
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
 - Estimation of Distribution
- 5 Conclusion

Stochastic Search

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

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

While not terminate

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

Everything depends on the definition of P and F_θ

deterministic algorithms are covered as well

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

Natural template for *Estimation of Distribution Algorithms*

Evolution Strategies and Normal Estimation of Distribution Algorithms

New search points are sampled normally distributed

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

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

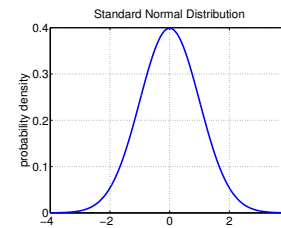
where

- 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

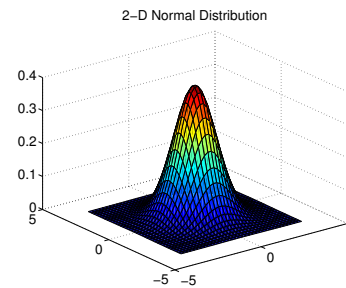
The question remains how to update m , \mathbf{C} , and σ .

Normal Distribution

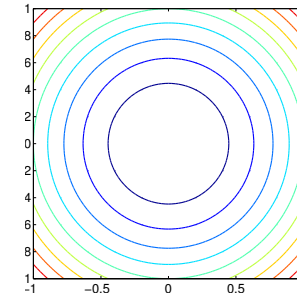
Isotropic Case



probability density of 1-D standard normal distribution



2-D

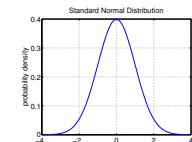


The Multi-Variate (n -Dimensional) Normal Distribution

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

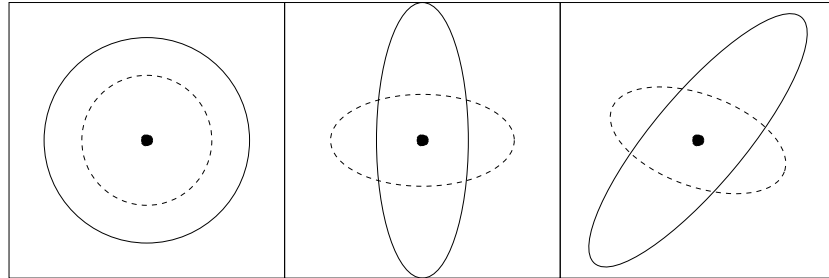
The **mean** value m

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



The **covariance matrix C** determines the shape. It has a valuable **geometrical interpretation**: any covariance matrix can be uniquely identified with the iso-density ellipsoid $\{x \in \mathbb{R}^n \mid x^T C^{-1} x = 1\}$

Lines of Equal Density



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

$\mathcal{N}(m, \mathbf{C}) \sim m + \mathbf{C}^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$ **(n² + n)/2 degrees of freedom** components are correlated

...CMA

Evolution Strategies

$(\mu \dagger \lambda)$ μ : # parents, λ : # offspring
 + selection in $\{\text{parents}\} \cup \{\text{offspring}\}$
 , selection in $\{\text{offspring}\}$

(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?

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

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

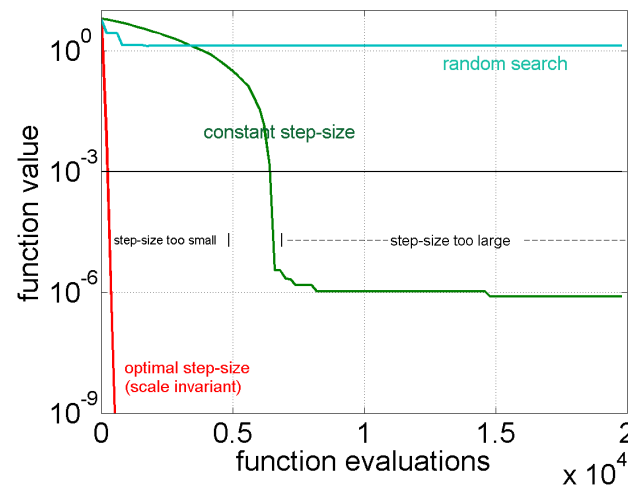
$$m \leftarrow \sum_{i=1}^{\mu} w_i x_{i:\lambda} = m + \sigma \underbrace{\sum_{i=1}^{\mu} w_i y_{i:\lambda}}_{=: y_w}$$

where

$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1$$

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

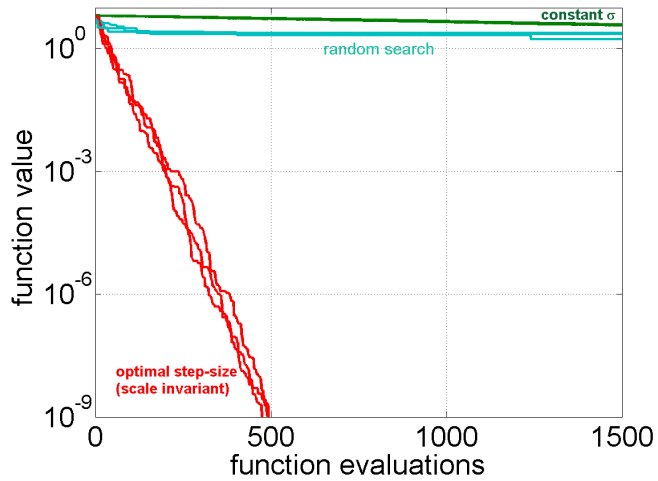
Why Step-Size Control?



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

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

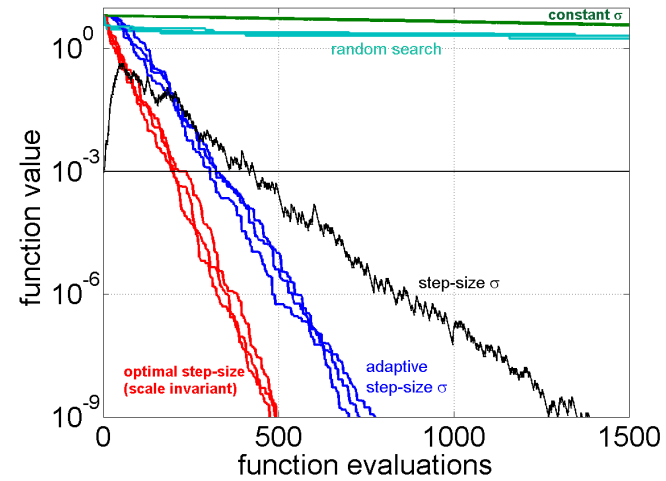
Why Step-Size Control?



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

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

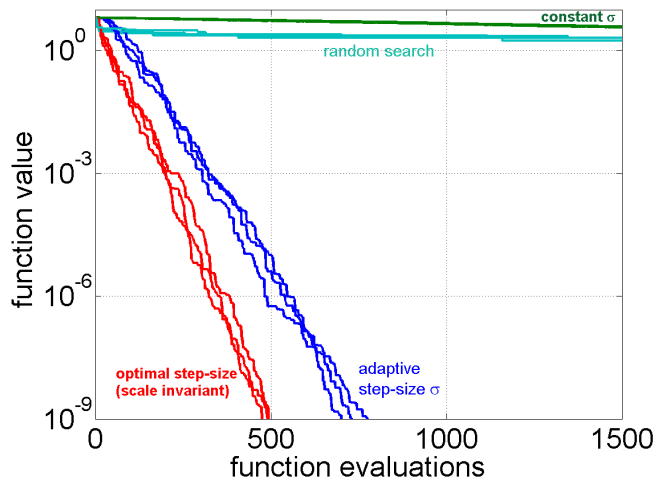
Why Step-Size Control?



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

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

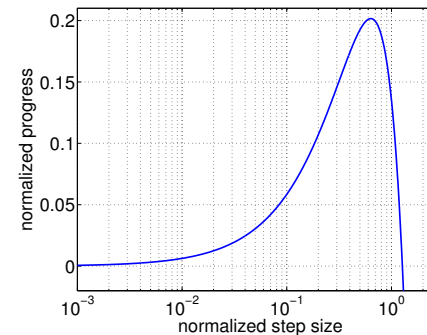
Why Step-Size Control?



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

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

Why Step-Size Control?



evolution window for the step-size on the sphere function

evolution window refers to the step-size interval where reasonable performance is observed

Methods for Step-Size Control

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

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

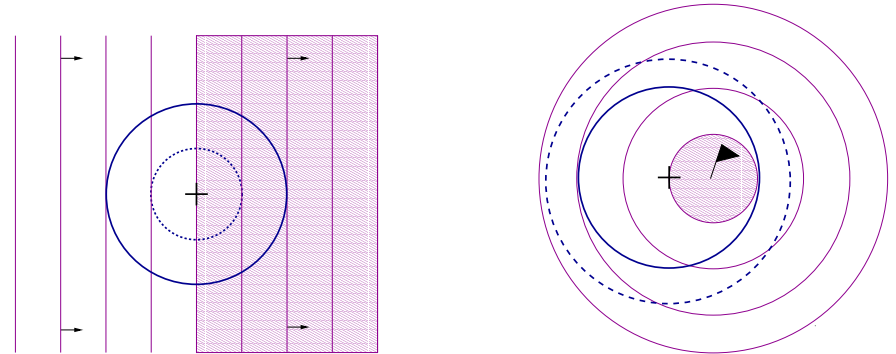
^bSchumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

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

^dHansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.* 9(2)

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

One-fifth success rule



Proba of success (p_s)

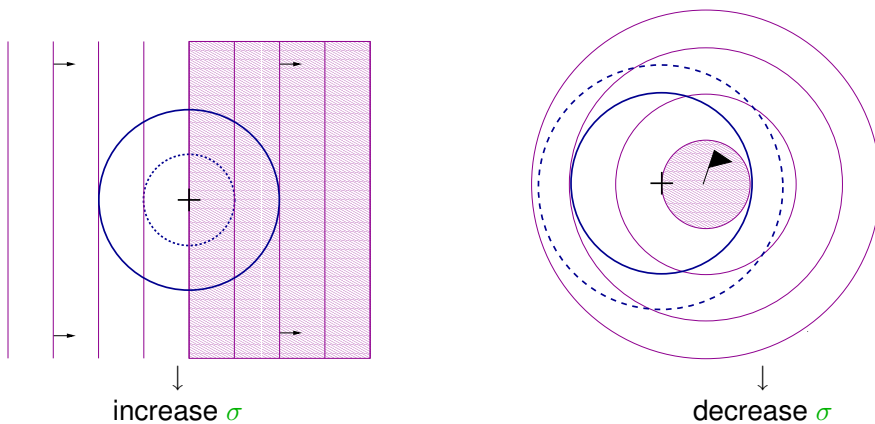
1/2

Proba of success (p_s)

1/5

“too small”

One-fifth success rule



increase σ

decrease σ

One-fifth success rule

Let p_s : # of successful offspring / generation

$$\sigma \leftarrow \sigma \times \exp\left(\frac{1}{3} \times \frac{p_s - p_{\text{target}}}{1 - p_{\text{target}}}\right)$$

Increase σ if $p_s > p_{\text{target}}$
 Decrease σ if $p_s < p_{\text{target}}$

(1 + 1)-ES

$$p_{\text{target}} = 1/5$$

IF offspring better parent

$$p_s = 1$$

ELSE

$$p_s = 0$$

Self-adaptation

in a $(1, \lambda)$ -ES

MUTATE for $i = 1, \dots, \lambda$

step-size
parent

$$\sigma_i \leftarrow \sigma \exp(\tau N(0, 1))$$

$$x_i \leftarrow x + \sigma_i \mathcal{N}(0, \mathbf{I})$$

EVALUATE

SELECT

Best offspring x_* with its step-size σ_*

Rationale

Unadapted step-size won't produce successive good individuals
 "The step-size are adjusted by the evolution itself"

Path Length Control

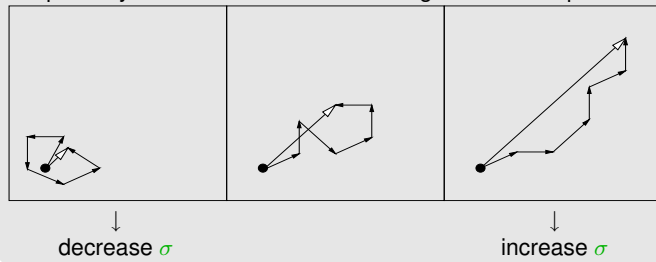
The Concept

$$x_i = m + \sigma y_i$$

$$m \leftarrow m + \sigma y_w$$

Measure the length of the evolution path

the pathway of the mean vector m in the generation sequence



loosely speaking steps are

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

Path Length Control

The Equations

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

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

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

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

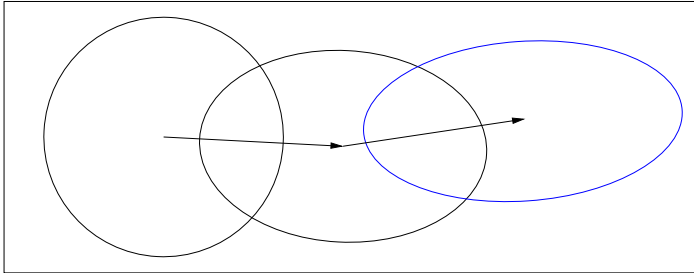
...CMA in a nutshell

- 1 Problem Statement
- 2 Evolution Strategies and EDAs
- 3 Step-Size Control
- 4 Covariance Matrix Adaptation
 - Covariance Matrix Rank-One Update
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank- μ Update
 - Estimation of Distribution
- 5 Conclusion

Covariance Matrix Adaptation

Rank-One Update

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



new distribution,

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

the ruling principle: the adaptation increases the probability of successful steps, \mathbf{y}_w , to appear again

... equations

Covariance Matrix Adaptation

Rank-One Update

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

While not terminate

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

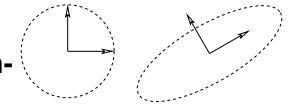
$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}$$

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

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

covariance matrix adaptation

- learns all **pairwise dependencies** between variables
off-diagonal entries in the covariance matrix reflect the dependencies
- conducts a **principle component analysis (PCA)** of steps \mathbf{y}_w , sequentially in time and space
eigenvectors of the covariance matrix \mathbf{C} are the principle components / the principle axes of the mutation ellipsoid
- learns a new, **rotated problem representation** and a **new metric** (Mahalanobis)
components are independent (only) in the new representation
- approximates the inverse Hessian on quadratic functions
overwhelming empirical evidence, proof is in progress



... cumulation, rank-μ, step-size control

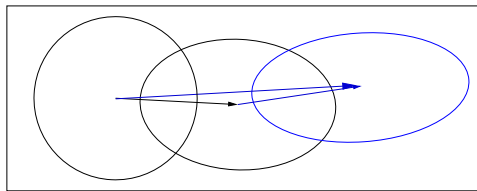
- 1 Problem Statement
- 2 Evolution Strategies and EDAs
- 3 Step-Size Control
- 4 Covariance Matrix Adaptation
 - Covariance Matrix Rank-One Update
 - Cumulation—the Evolution Path
 - Covariance Matrix Rank-μ Update
 - Estimation of Distribution
- 5 Conclusion

Cumulation

The Evolution Path

Evolution Path

Conceptually, the evolution path is the **path** the strategy takes **over a number of generation steps**. It can be expressed as a sum of consecutive **steps** of the mean ***m***.



An exponentially weighted sum of steps y_w is used

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

The recursive construction of the evolution path (cumulation):

$$p_c \leftarrow \underbrace{(1 - c_c)}_{\text{decay factor}} p_c + \underbrace{\sqrt{1 - (1 - c_c)^2}}_{\text{normalization factor}} \underbrace{y_w}_{\text{input, } \frac{m - m_{\text{old}}}{\sigma}}$$

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

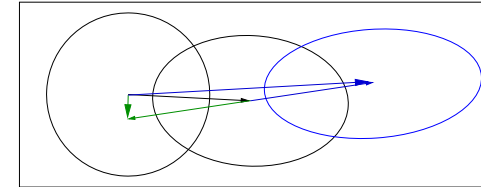
- *exponential smoothing* in time series, forecasting
- exponentially weighted *moving average*
- *iterate averaging* in stochastic approximation
- *momentum* in the back-propagation algorithm for ANNs
- ...

... why?

Cumulation

Utilizing the Evolution Path

We used $y_w y_w^T$ for updating **C**. Because $y_w y_w^T = -y_w (-y_w)^T$ the sign of y_w is neglected. The sign information is (re-)introduced by using the *evolution path*.



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

where $\mu_w = \frac{1}{\sum w_i^2}$, $c_c \ll 1$.

... equations

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 $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$** .^a

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

... rank-1 update

Rank- μ Update

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

The rank- μ update extends the update rule for **large population sizes** λ using $\mu > 1$ vectors to update \mathbf{C} at each generation step.

The matrix

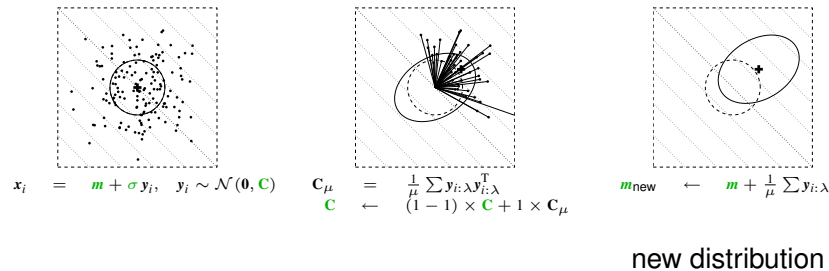
$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

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

The rank- μ update then reads

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

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



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

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

The rank- μ update

- increases the possible learning rate in large populations roughly from $2/n^2$ to μ_w/n^2
- can reduce the number of necessary **generations** roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)^3$ given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ 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- μ update can be combined. . .

³Hansen, 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

Estimation of Distribution Algorithms

- Estimate a distribution that (re-)samples the parental population.
- All parameters of the distribution θ are estimated from the given population.

Example: EMNA (Estimation of Multi-variate Normal Algorithm)

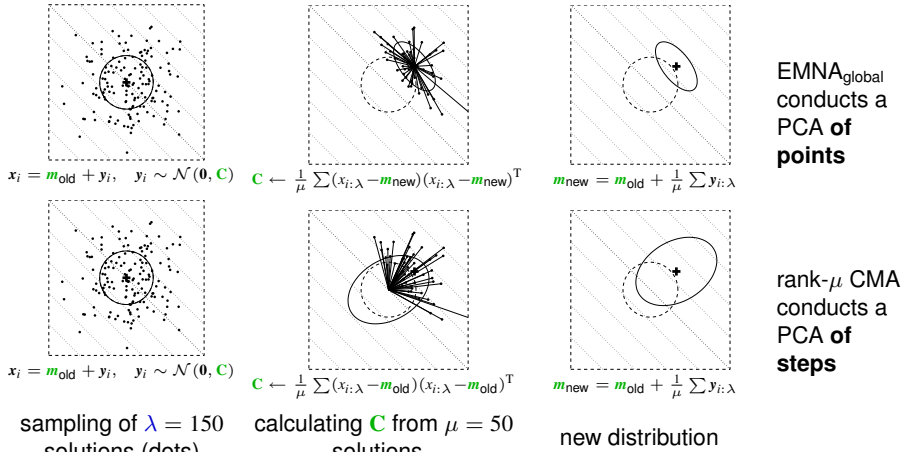
Initialize $\mathbf{m} \in \mathbb{R}^n$, and $\mathbf{C} = \mathbf{I}$

While not terminate

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

Larrañaga and Lozano 2002. *Estimation of Distribution Algorithms*

Estimation of Multivariate Normal Algorithm EMNA_{global} versus rank- μ CMA⁴



The CMA-update yields a larger variance in particular in gradient direction, because m_{new} is the minimizer for the variances when calculating C

⁴ 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

Conclusion

- 1 Problem Statement
- 2 Evolution Strategies and EDAs
- 3 Step-Size Control
- 4 Covariance Matrix Adaptation
- 5 Conclusion

What did we achieve?

- 1 Covariance matrix adaptation: reduce any convex quadratic function

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

to the sphere model

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

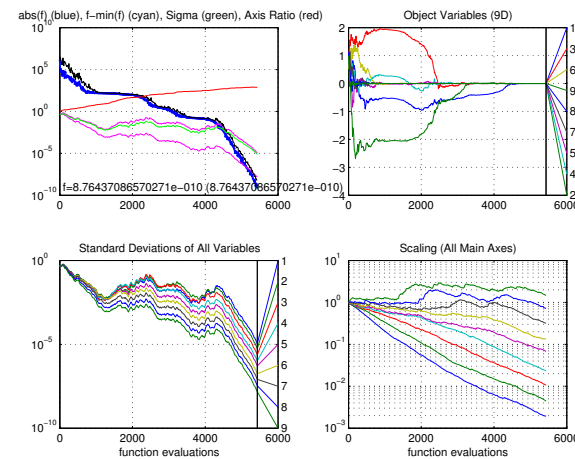
lines of equal density align with lines of equal fitness $C \propto H^{-1}$ without use of derivatives

- 2 Step-size control: converge log-linearly on the sphere
- 3 Rank-based selection: the same holds for any $g(f(x)) = g(x^T H x)$
 $g: \mathbb{R} \rightarrow \mathbb{R}$ strictly monotonic (order preserving)

Conclusion

Experimentum Crucis (1)

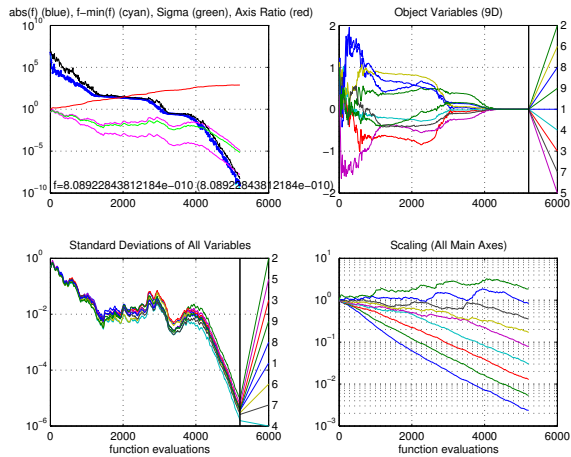
f convex quadratic, separable



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

Experimentum Crucis (2)

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



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

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

... internal parameters

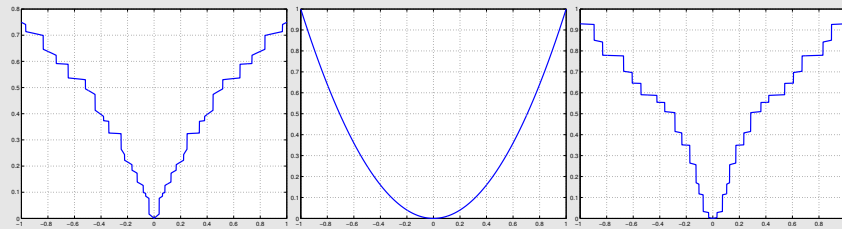
Invariance Under Strictly Monotonically Increasing Functions

Rank-based algorithms

Selection based on the rank:

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

Update of all parameters uses only the rank



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