

Advanced Optimization

Lecture 3: Randomized Algorithms for Continuous Problems

Master AIC

Université Paris-Saclay, Orsay, France

Anne Auger

INRIA Saclay – Ile-de-France



Dimo Brockhoff

INRIA Saclay – Ile-de-France

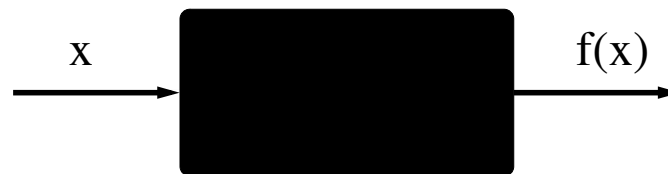
Problem Statement

Continuous Domain Search/Optimization

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

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

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

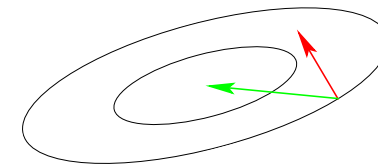
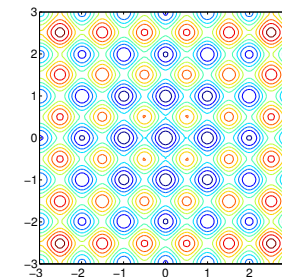
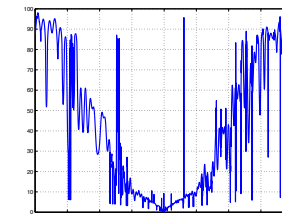
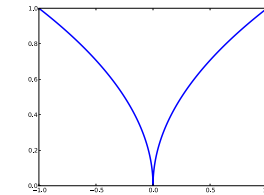


- ▶ gradients are not available or not useful
- ▶ problem domain specific knowledge is used only within the black box, e.g. within an appropriate encoding
- ▶ Search **costs**: number of function evaluations

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

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 100 points onto a real interval, say $[0, 1]$. To get **similar coverage**, in terms of distance between adjacent points, of the 10-dimensional space $[0, 1]^{10}$ would require $100^{10} = 10^{20}$ points. A 100 points appear now as isolated points in a vast empty space.

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

Separable Problems

Definition (Separable Problem)

A function f is separable if

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

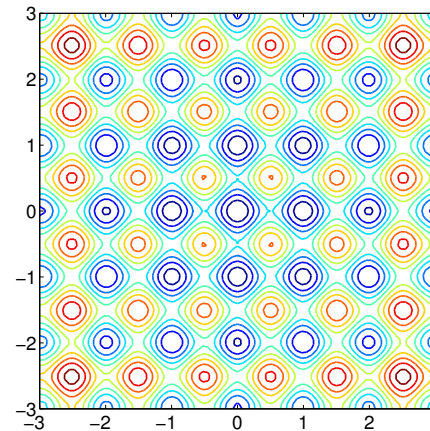
⇒ it follows that f can be optimized in a sequence of n independent 1-D optimization processes

Example: Additively decomposable functions

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

Rastrigin function

$$f(\mathbf{x}) = 10n + \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i))$$



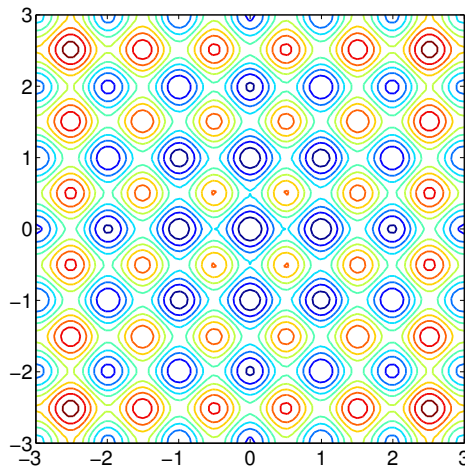
Non-Separable Problems

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

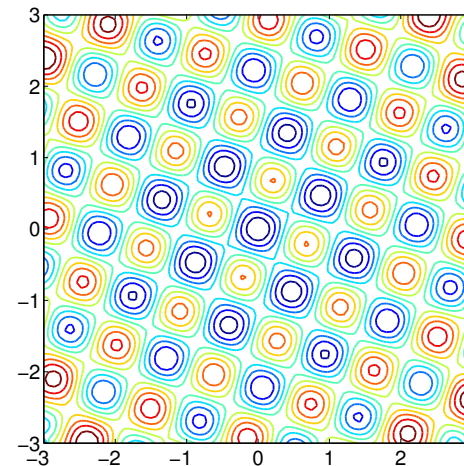
Rotating the coordinate system

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

R rotation matrix



R
→



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

Ill-Conditioned Problems

- ▶ If f is convex quadratic, $f : \mathbf{x} \mapsto \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} = \frac{1}{2} \sum_i h_{i,i} x_i^2 + \frac{1}{2} \sum_{i \neq j} h_{i,j} x_i x_j$, with \mathbf{H} positive, definite, symmetric matrix

\mathbf{H} is the Hessian matrix of f

- ▶ ill-conditioned means a high condition number of Hessian Matrix \mathbf{H}

$$\text{cond}(\mathbf{H}) = \frac{\lambda_{\max}(\mathbf{H})}{\lambda_{\min}(\mathbf{H})}$$

Example / exercise

The **level-sets** of a function are defined as

$$\mathcal{L}_c = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) = c\}, c \in \mathbb{R}.$$

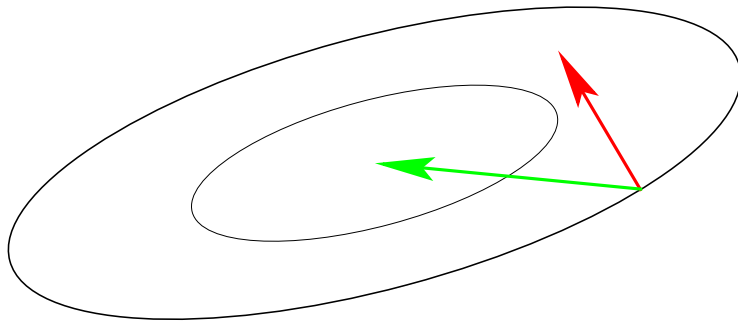
Consider the objective function $f(\mathbf{x}) = \frac{1}{2}(x_1^2 + 9x_2^2)$

1. Plot the level sets of f .
2. Compute the condition number of the Hessian matrix of f , relate it to the axis ratio of the level sets of f .
3. Generalize 1. and 2. to a general convex-quadratic function.

Ill-conditioned Problems

consider the curvature of the level sets of a function

ill-conditioned means “squeezed” lines of equal function value (high curvatures)



gradient direction $-f'(x)^T$

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

Condition number equals nine here. Condition numbers up to 10^{10} are not unusual in real world problems.

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

Everything depends on the definition of P and F_θ

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*

A Simple Example: The Pure Random Search

Also an Ineffective Example

The Pure Random Search

- ▶ Sample **uniformly** at random a solution
- ▶ Return the best solution ever found

Exercise

See the exercise on the document "Exercices - class 1".

Non-adaptive Algorithm

For the pure random search $P(x|\theta)$ is independent of θ (i.e. no θ to be adapted): the algorithm is "blind"

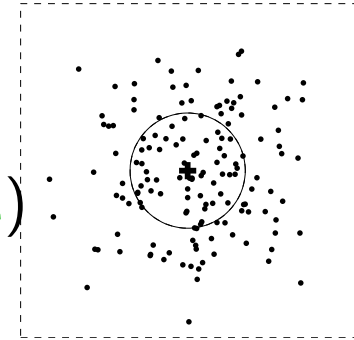
In this class: present algorithms that are "much better" than that

Evolution Strategies

New search points are sampled normally distributed

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

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

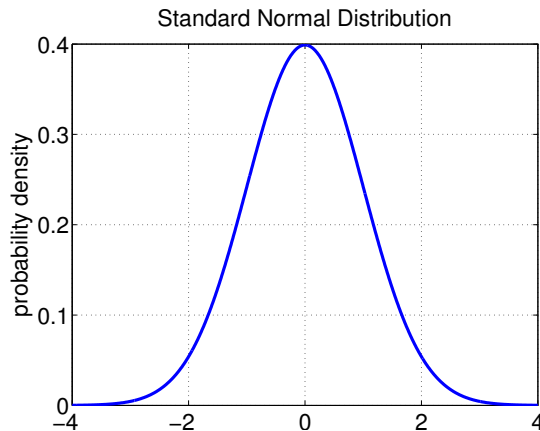
- ▶ the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- ▶ the so-called **step-size** $\sigma \in \mathbb{R}_+$ controls the *step length*
- ▶ the **covariance matrix** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

here, all new points are sampled with the same parameters

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

Normal Distribution

1-D case



probability density of the 1-D standard normal distribution $\mathcal{N}(0, 1)$

(expected (mean) value, variance) = (0,1)

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

General case

- ▶ Normal distribution $\mathcal{N}(m, \sigma^2)$

(expected value, variance) = (m, σ^2)

density: $p_{m,\sigma}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right)$

- ▶ A normal distribution is entirely determined by its mean value and variance
- ▶ The family of normal distributions is closed under linear transformations: if X is normally distributed then a linear transformation $aX + b$ is also normally distributed
- ▶ **Exercise:** Show that $m + \sigma\mathcal{N}(0, 1) = \mathcal{N}(m, \sigma^2)$

Normal Distribution

General case

A random variable following a 1-D normal distribution is determined by its **mean value** m and **variance** σ^2 .

In the n -dimensional case it is determined by its **mean vector** and **covariance matrix**

Covariance Matrix

If the entries in a vector $\mathbf{X} = (X_1, \dots, X_n)^T$ are random variables, each with finite variance, then the covariance matrix Σ is the matrix whose (i, j) entries are the covariance of (X_i, X_j)

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = \mathbb{E} [(X_i - \mu_i)(X_j - \mu_j)]$$

where $\mu_i = \mathbb{E}(X_i)$. Considering the expectation of a matrix as the expectation of each entry, we have

$$\Sigma = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$$

Σ is symmetric, positive definite

The Multi-Variate (n -Dimensional) Normal Distribution

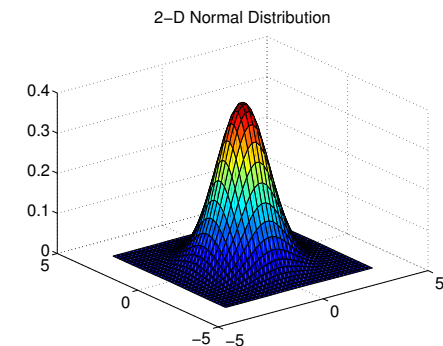
Any multi-variate normal distribution $\mathcal{N}(\mathbf{m}, \mathbf{C})$ is uniquely determined by its mean value $\mathbf{m} \in \mathbb{R}^n$ and its symmetric positive definite $n \times n$ covariance matrix \mathbf{C} .

$$\text{density: } p_{\mathcal{N}(\mathbf{m}, \mathbf{C})}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m})\right),$$

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

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

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

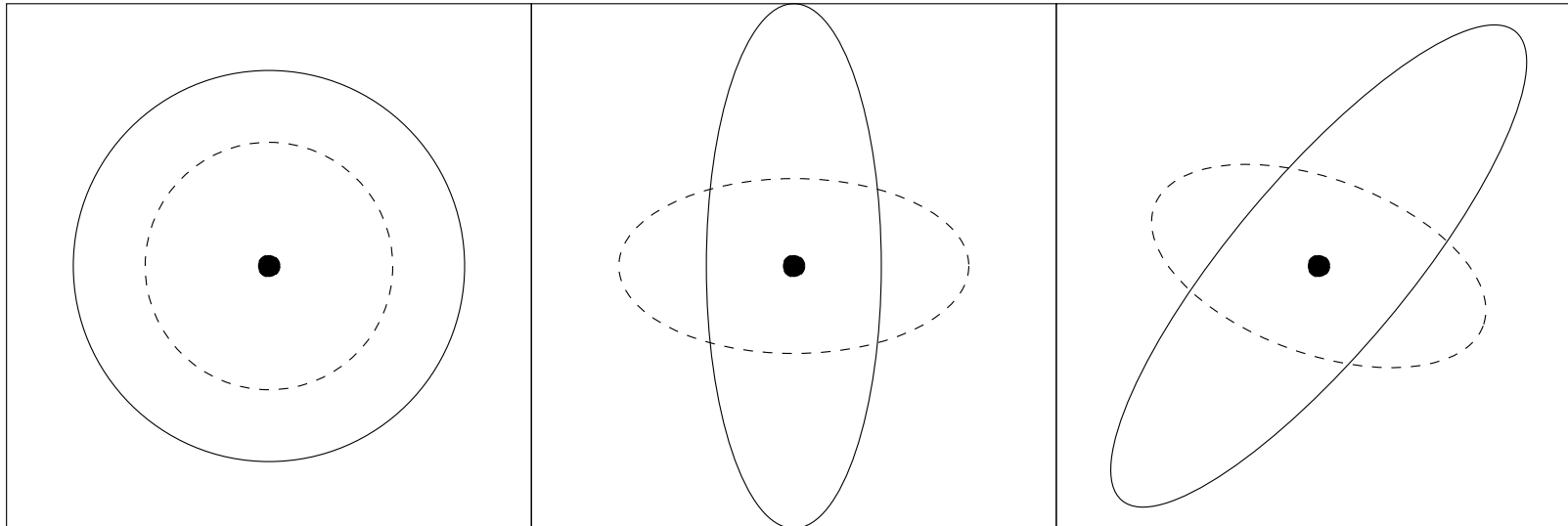


The **covariance matrix** \mathbf{C}

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

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

Lines of Equal Density



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

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

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

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

Adapting the mean ...

Evolution Strategies (ES)

Simple Update for Mean Vector

Let μ : # parents, λ : # offspring

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

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

(μ, λ) -ES: selection in $\{\text{offspring}\}$

ES algorithms emerged in the community of bio-inspired methods where a parallel between optimization and evolution of species as described by Darwin served in the origin as inspiration for the methods. Nowadays this parallel is mainly visible through the terminology used: candidate solutions are parents or offspring, the objective function is a fitness function, ...

$(1 + 1)$ -ES

Sample one offspring from parent m

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

If x better than m select

$$m \leftarrow x$$

The $(\mu/\mu, \lambda)$ -ES - Update of the mean vector

Non-elitist selection and intermediate (weighted) recombination

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

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

Notation: we denote $\mathbf{y}_{i:\lambda}$ the vector such that $\mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \mathbf{y}_{i:\lambda}$

Exercise: realize that $\mathbf{y}_{i:\lambda}$ is generally not distributed as $\mathcal{N}(\mathbf{0}, \mathbf{C})$

The new mean reads

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

where

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

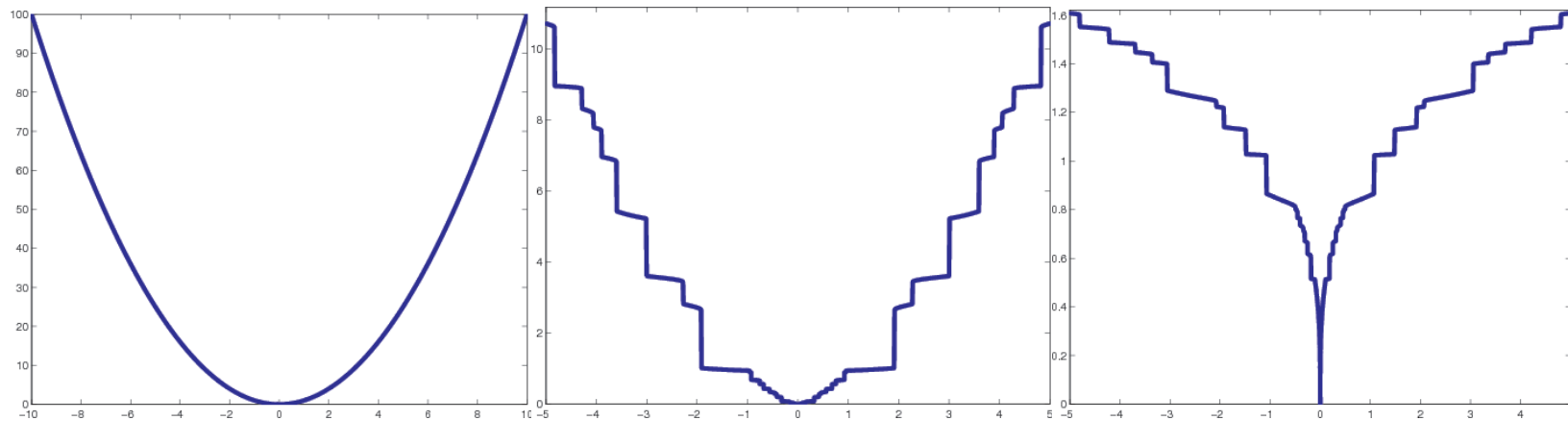
The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

Invariance Under Monotonically Increasing Functions

Rank-based algorithms

Update of all parameters uses only the ranks

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

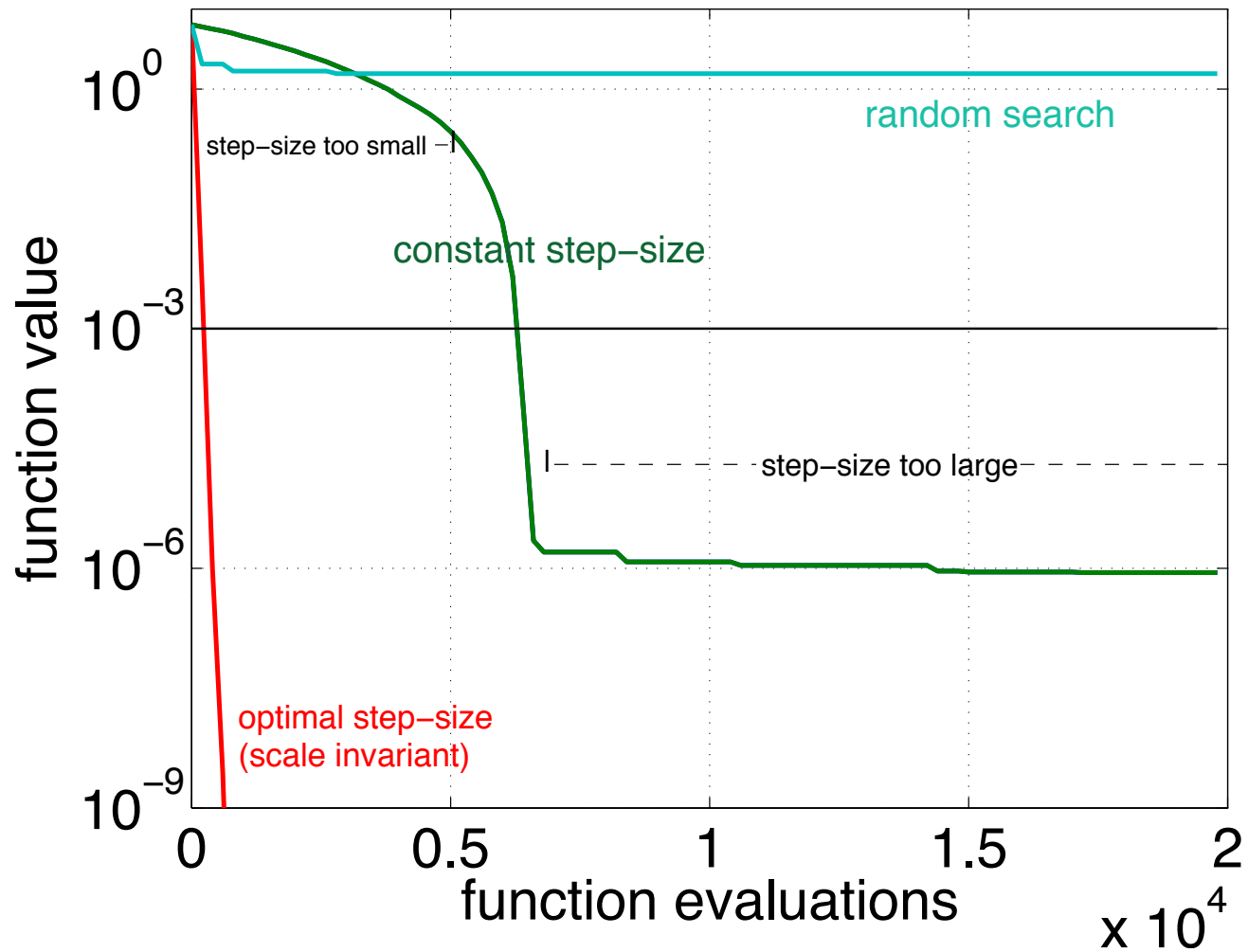


$$g(f(x_{1:\lambda})) \leq g(f(x_{2:\lambda})) \leq \dots \leq g(f(x_{\lambda:\lambda})) \quad \forall g$$

g is strictly monotonically increasing
g preserves ranks

Adapting the step-size ...

Why Step-Size Control?



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

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

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

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, *Evolutionstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog

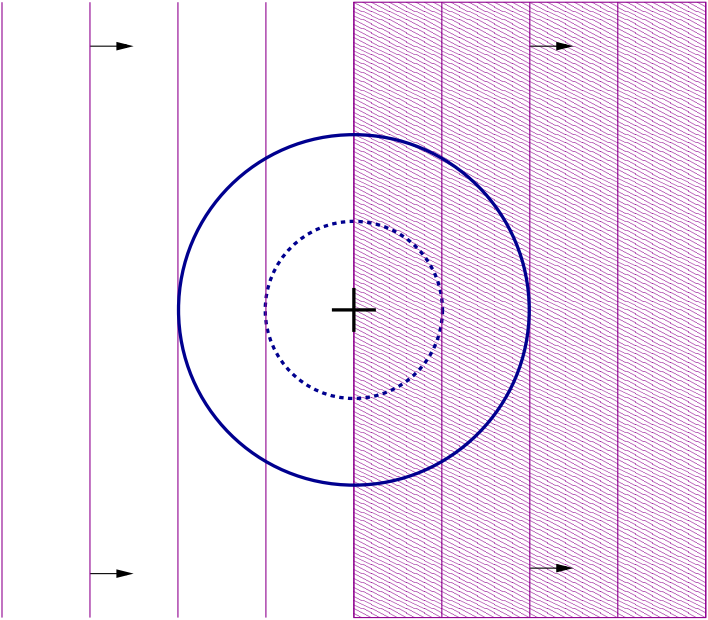
^bSchumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

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

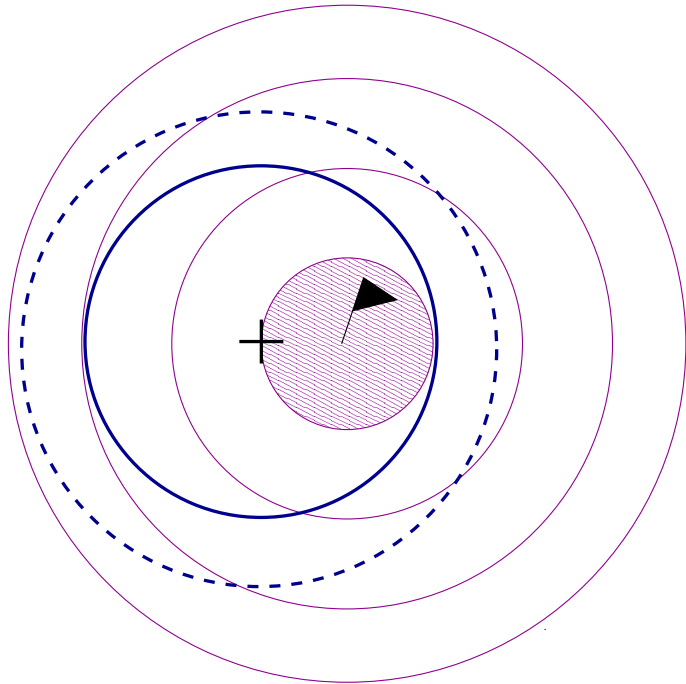
^dHansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.* 9(2)

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

One-fifth success rule

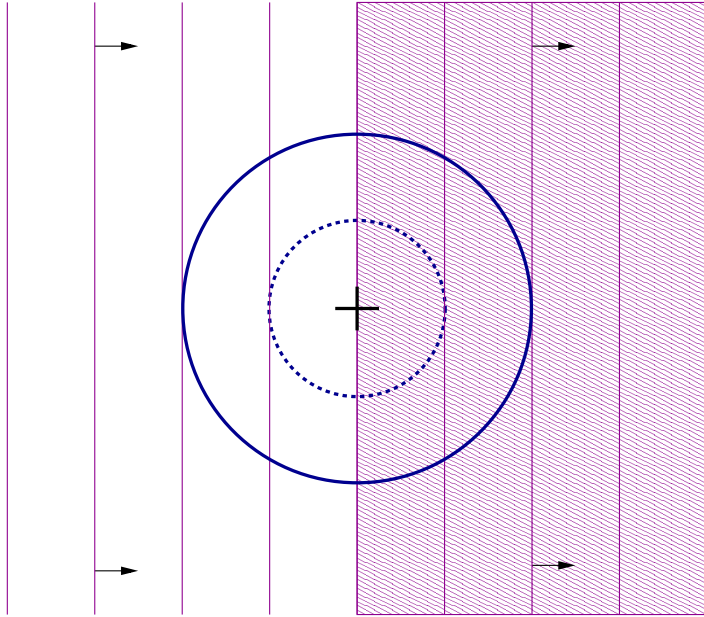


↓
increase σ



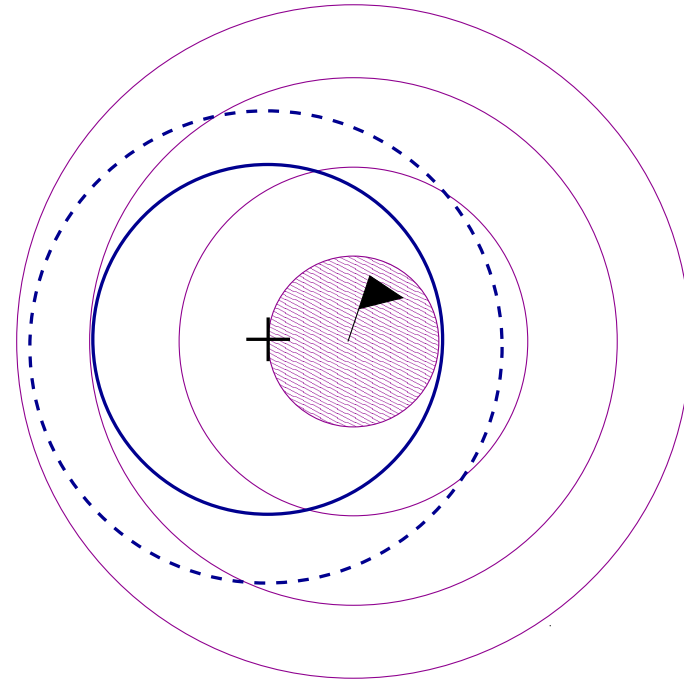
↓
decrease σ

One-fifth success rule



Probability of success (p_s)

$1/2$



Probability of success (p_s)

“too small”

One-fifth success rule

p_s : # of successful offspring / # offspring (per iteration)

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

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*

$$p_s = 1, \sigma \leftarrow \sigma \times \exp(1/3)$$

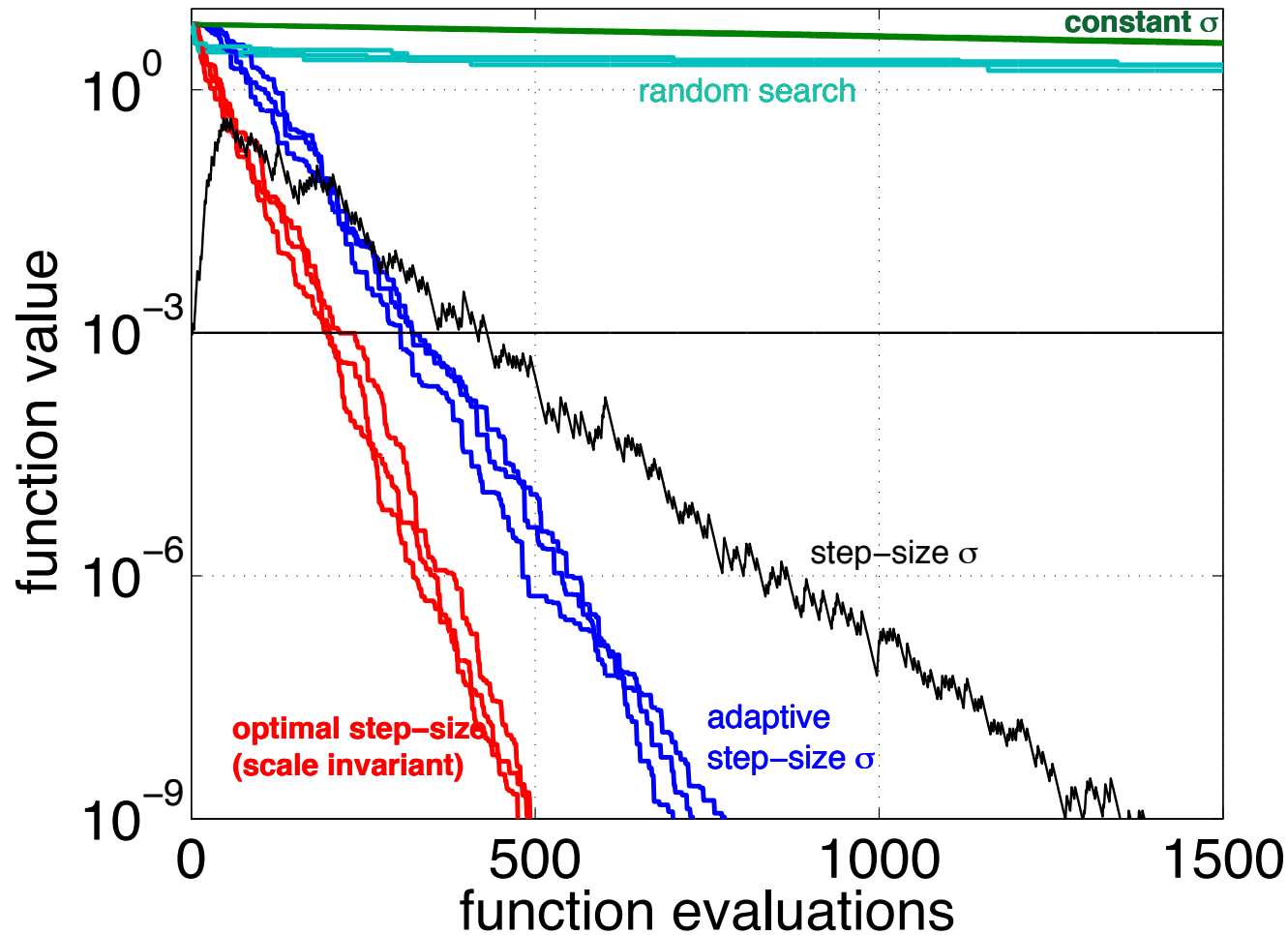
ELSE

$$p_s = 0, \sigma \leftarrow \sigma / \exp(1/3)^{1/4}$$

Step-size adaptation

What is achieved

(1 + 1)-ES with one-fifth success rule (blue)

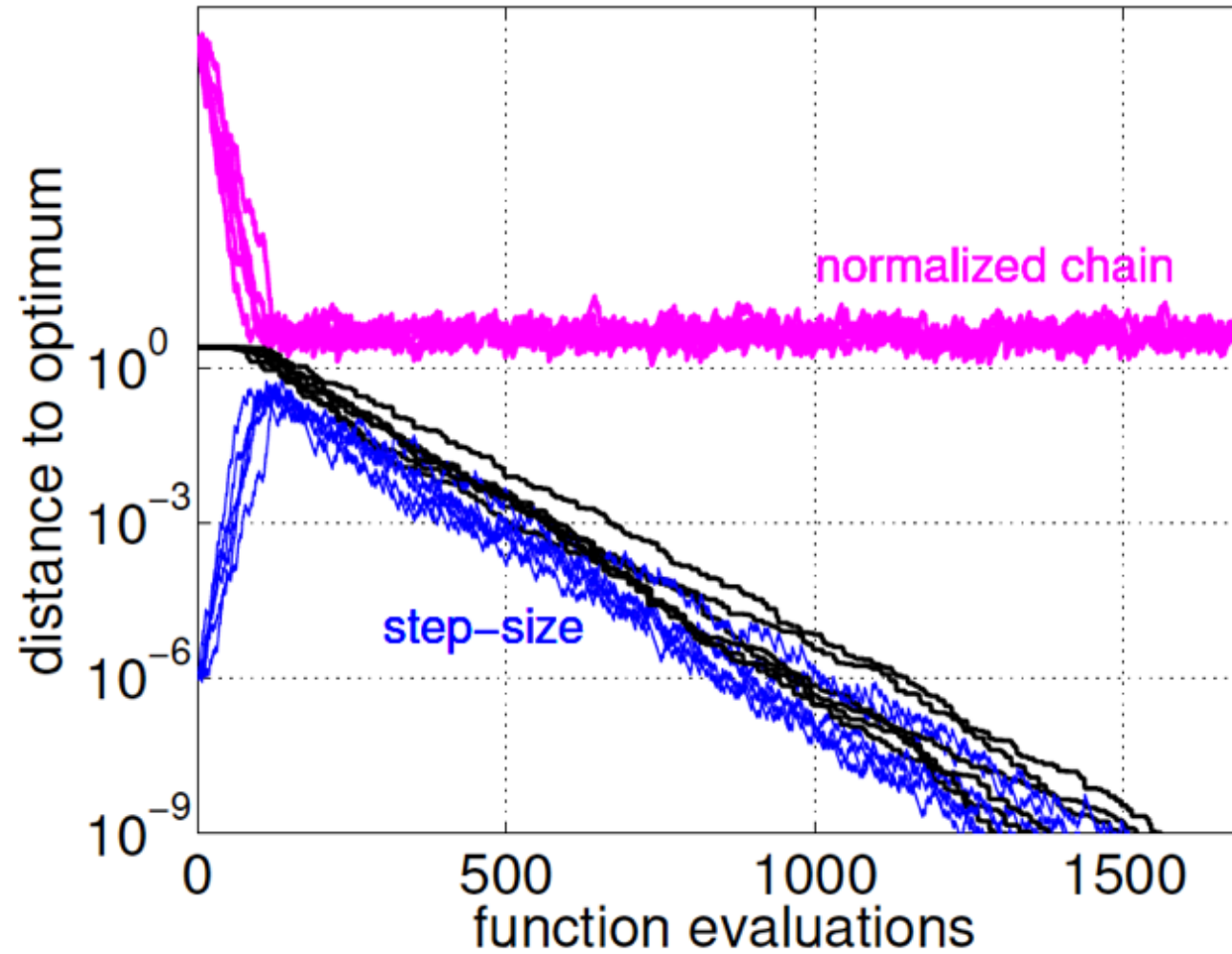


$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

Linear convergence

What do we achieve?

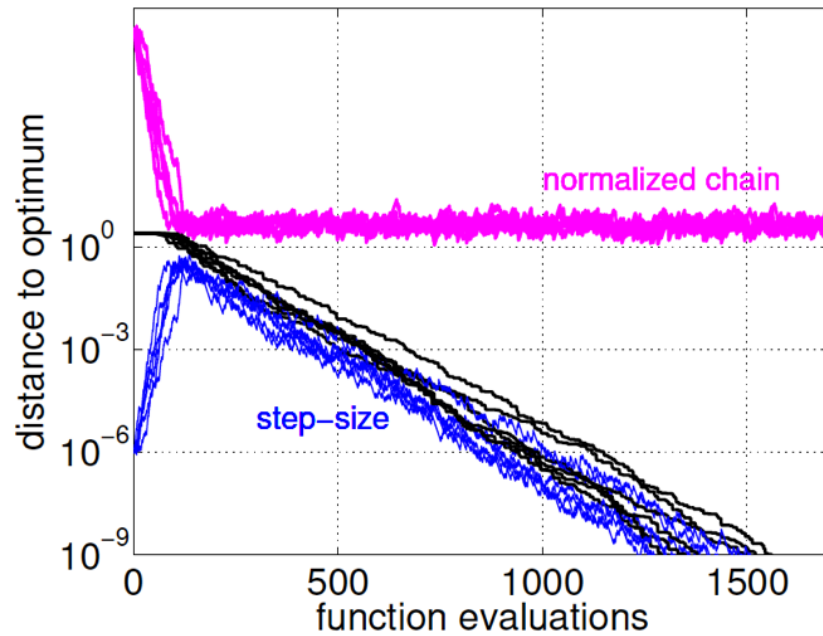


(1+1)-ES

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

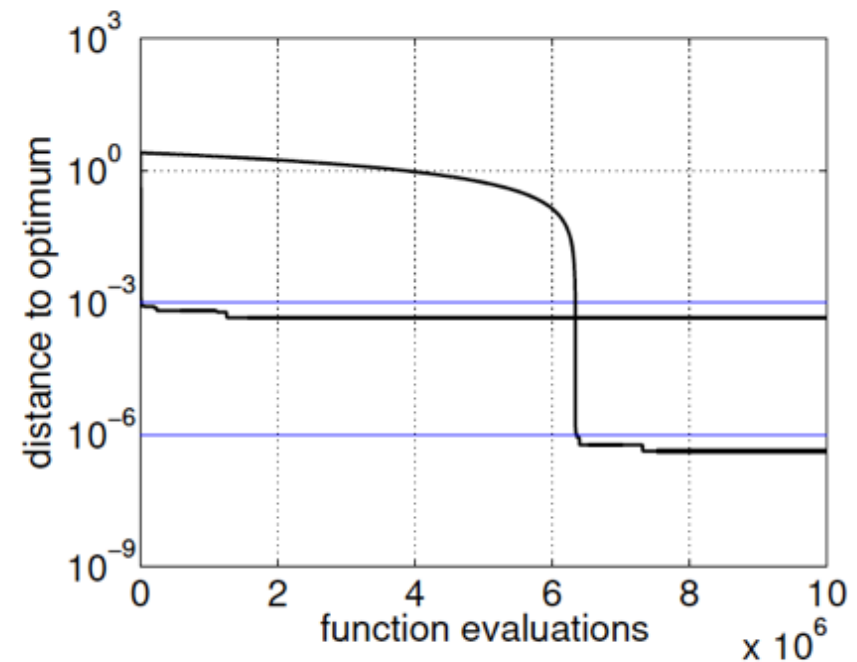
for $n = 10$

Adaptive versus Constant Step-size



adaptive step-size

constant step-size



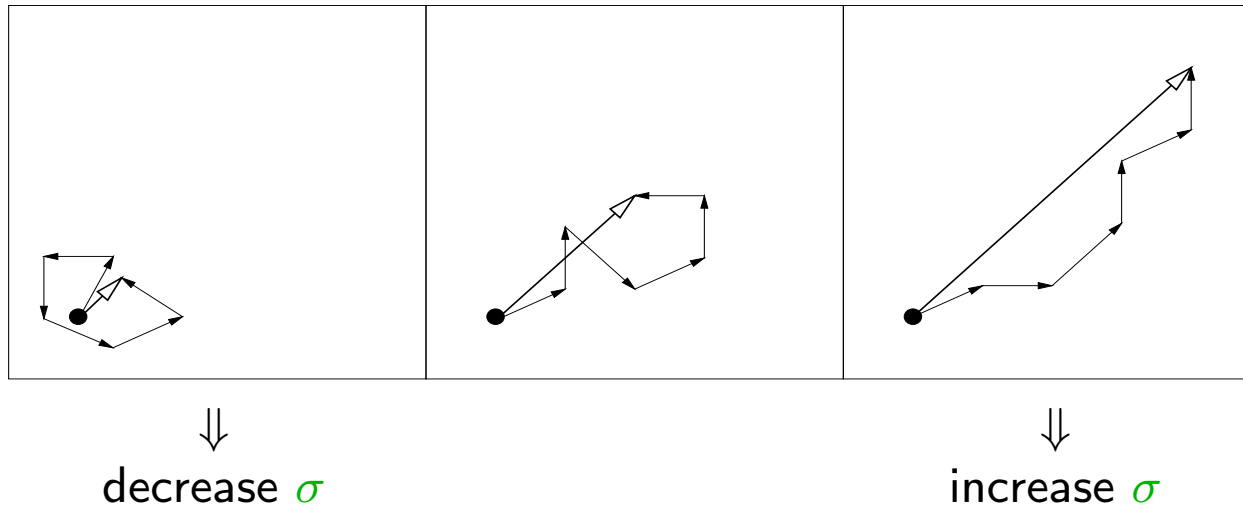
Path Length Control (CSA)

The Concept of Cumulative Step-Size Adaptation

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

Measure the length of the *evolution path*

the pathway of the mean vector \mathbf{m} in the iteration sequence



Path Length Control (CSA)

The Equations

Sampling of solutions, notations as on slide “The $(\mu/\mu, \lambda)$ -ES - Update of the mean vector” with \mathbf{C} equal to the identity.

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

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

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

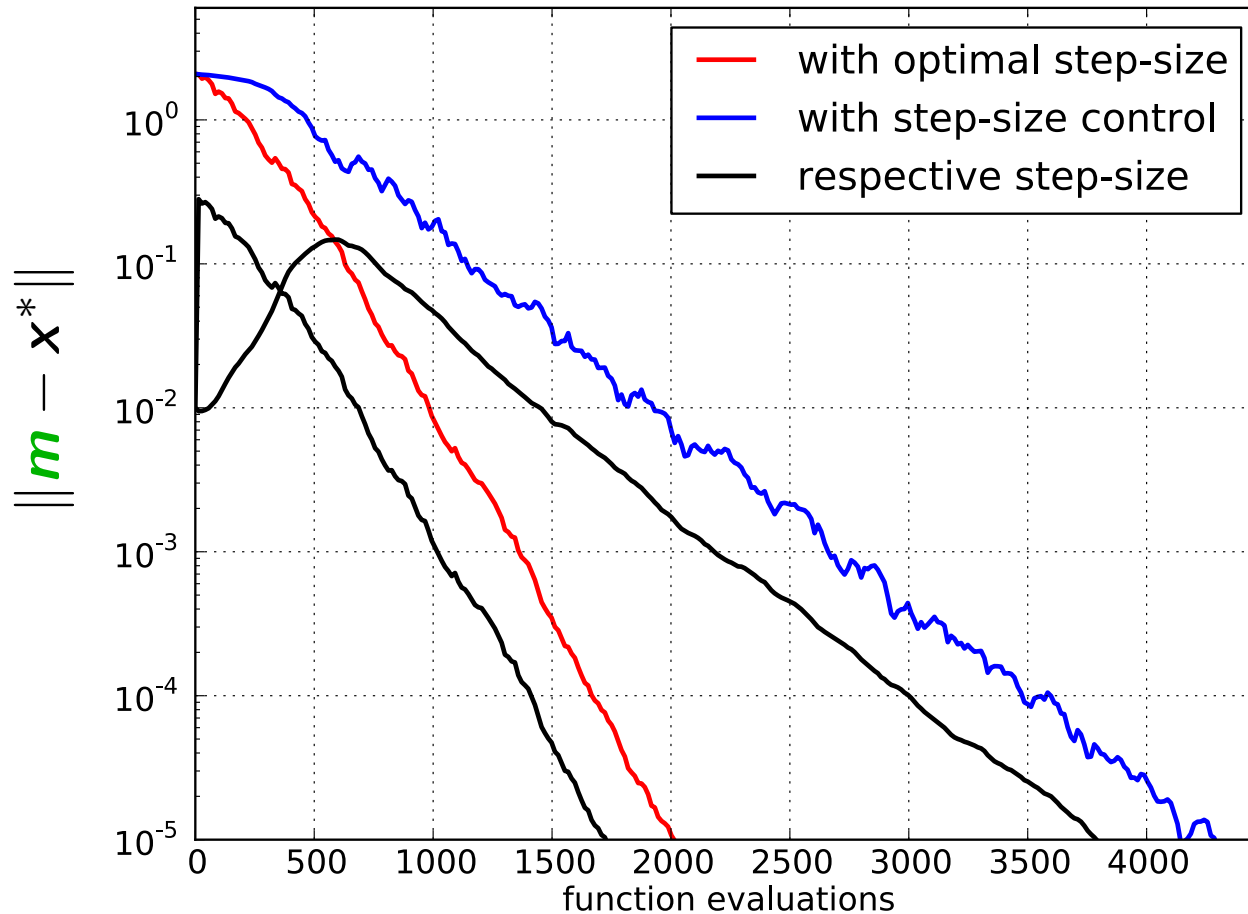
$$\sigma \leftarrow \sigma \times \underbrace{\exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)}_{\text{update step-size}}$$

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

Step-size adaptation

What is achieved

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

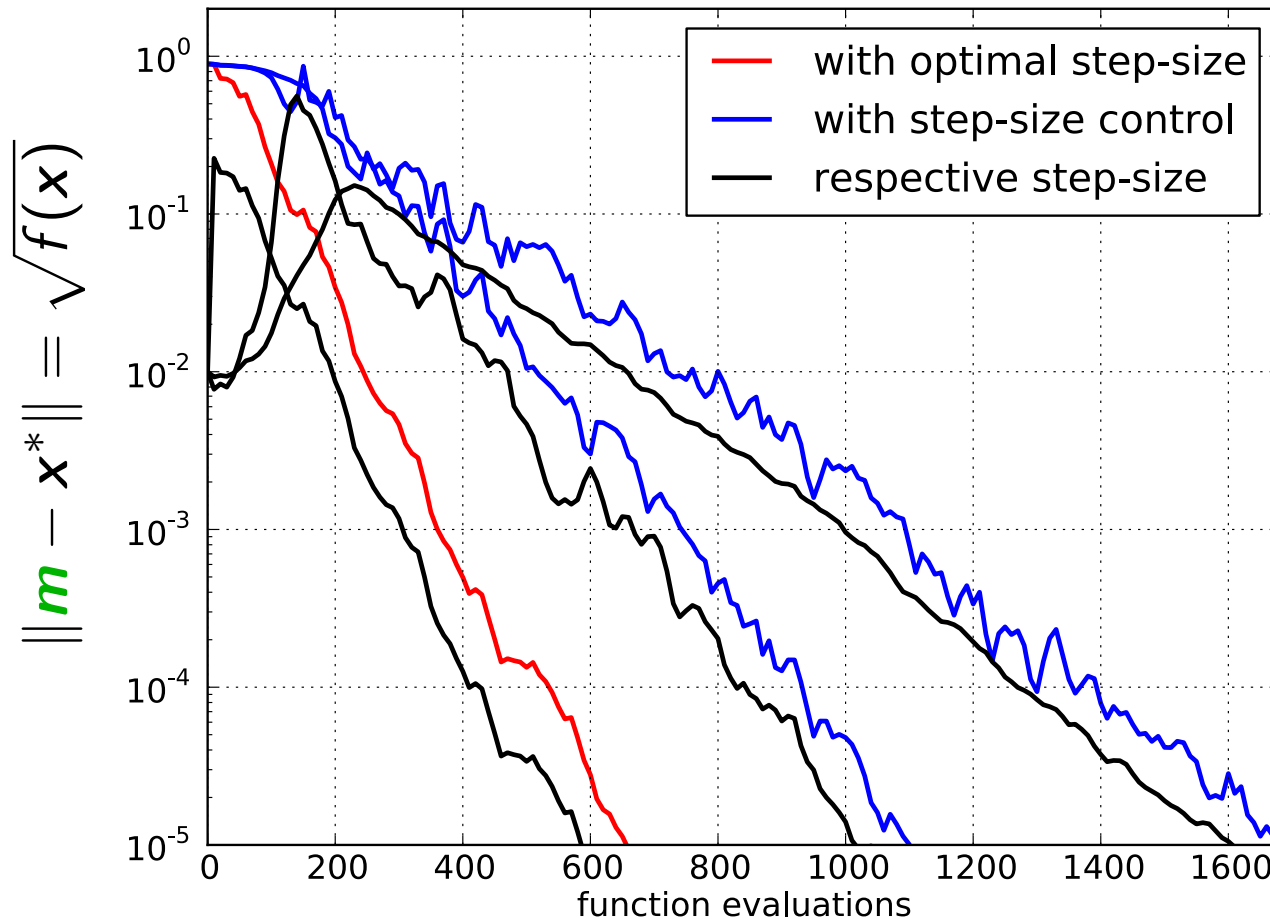


$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

Why Step-Size Control?

(5/5_w, 10)-ES



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

for $n = 10$

and

$$\mathbf{x}^0 \in [-0.2, 0.8]^n$$

comparing optimal versus default damping parameter d_σ :

$$\frac{1700}{1100} \approx 1.5$$

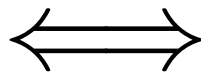
On linear convergence ...

Hitting Time versus Convergence

Finite hitting time for all epsilon

$$T_\epsilon = \inf \{t \in \mathbb{N}, \mathbf{X}_t \in B(\mathbf{x}^*, \epsilon)\}$$

$$T_\epsilon < \infty \text{ for all } \epsilon > 0$$



under some regularity conditions on
the algorithm and the function
e.g.) (1+1)-ES on a spherical function

Convergence towards the optimum

$$\lim_{t \rightarrow \infty} \mathbf{X}_t = \mathbf{x}^*$$

$$\iff \forall \epsilon > 0, \exists T_\epsilon < \infty \text{ such that } \|\mathbf{X}_t - \mathbf{x}^*\| < \epsilon \text{ for all } t \geq T_\epsilon$$

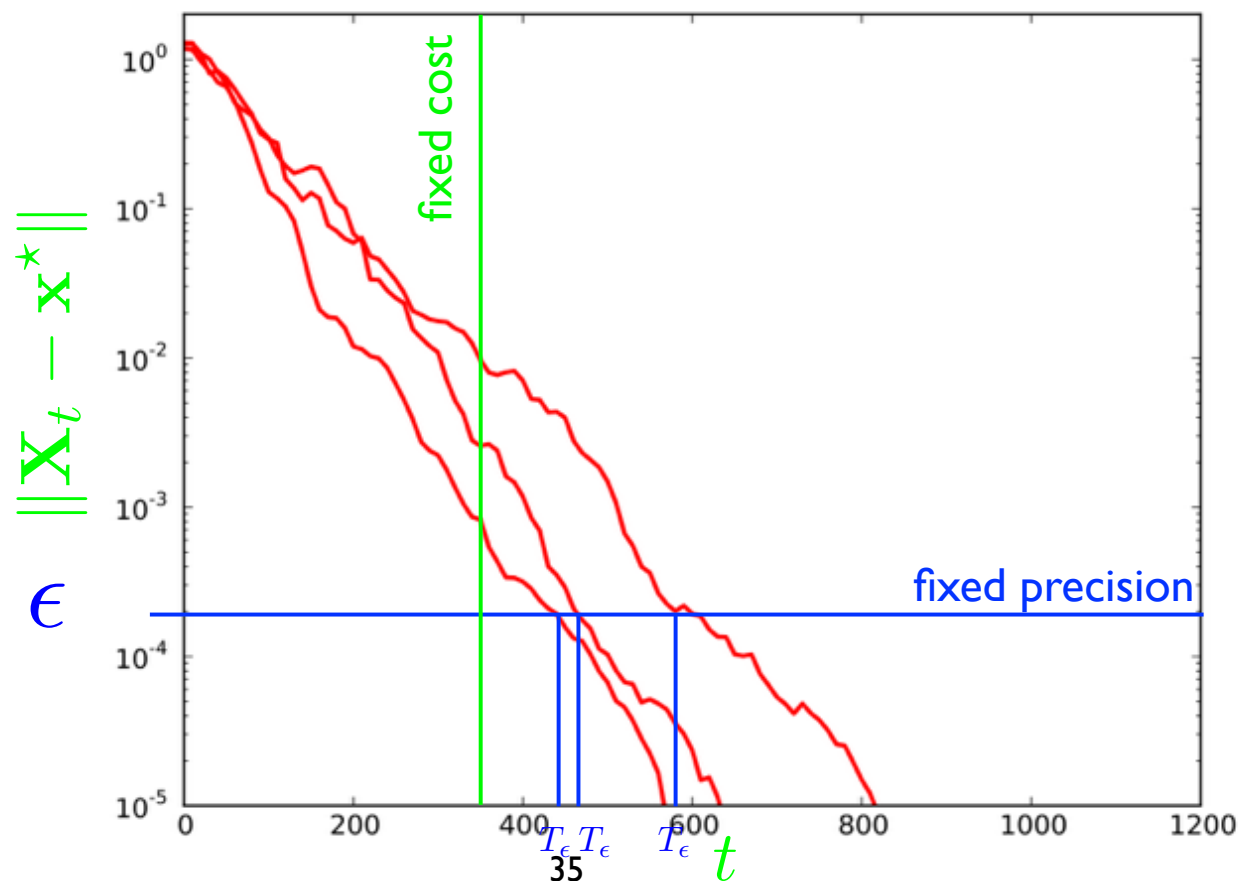
translate that an algorithm approximates the
optimum with **arbitrary** precision

Hitting Time versus Convergence

two side of a coin, measuring

the hitting time T_ϵ given a fixed precision ϵ

the precision $\|\mathbf{X}_t - \mathbf{x}^*\|$ (or ϵ) given the iteration number t



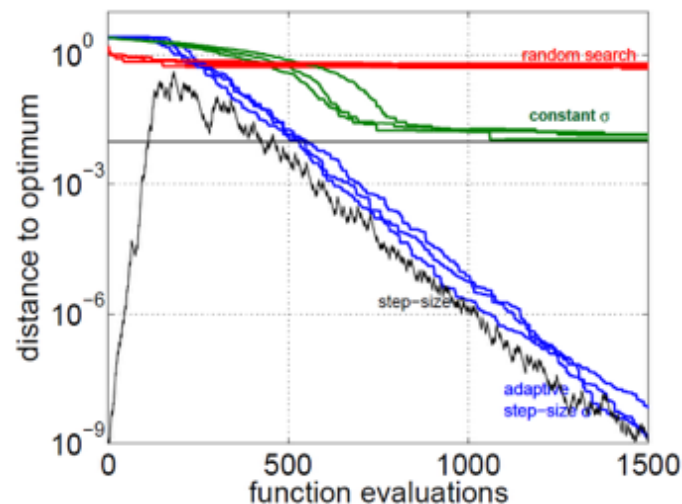
On Convergence alone ...

A theoretical convergence result is a “guarantee” that the algorithm will approach the solution in **infinite** time

$$\lim_{t \rightarrow \infty} \mathbf{X}_t = \mathbf{x}^*$$

often the first/only question investigated about an optimization algorithm

But a convergence result alone is pretty meaningless **in practice** as it does not tell how fast the algorithm converges



need to quantify how fast the optimum is approached

Quantifying How Fast the Optimum is Approached

For a fixed dimension

convergence speed of
 \mathbf{X}_t towards \mathbf{x}^*



dependency in ϵ of T_ϵ
find $\epsilon \mapsto \tau(\epsilon, n)$

Scaling wrt the dimension

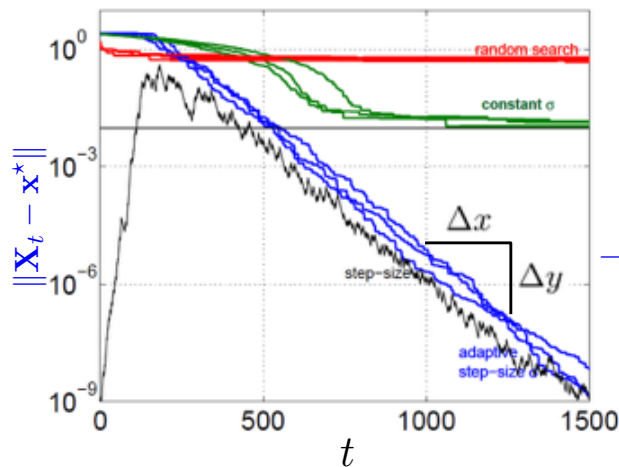
dependency of convergence
rate wrt n



find $n \mapsto \tau(\epsilon, n)$

Compromises to obtain such results:
asymptotic in n , in ϵ / t

Linear Convergence



$$\frac{\|\mathbf{X}_{t+1} - \mathbf{x}^*\|}{\|\mathbf{X}_t - \mathbf{x}^*\|} \approx \exp\left(-\frac{c}{n}\right)$$

$$\log \frac{\|\mathbf{X}_{t+1} - \mathbf{x}^*\|}{\|\mathbf{X}_t - \mathbf{x}^*\|} \approx -\frac{c}{n}$$

$$\log \frac{\|\mathbf{X}_t - \mathbf{x}^*\|}{\|\mathbf{X}_0 - \mathbf{x}^*\|} \approx -\frac{c}{n}t$$

Different formal statements (not exactly equivalent)

almost surely

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|\mathbf{X}_t - \mathbf{x}^*\|}{\|\mathbf{X}_0 - \mathbf{x}^*\|} = -\frac{c}{n}$$

in expectation

$$\frac{\mathbb{E} [\|\mathbf{X}_{t+1} - \mathbf{x}^*\|]}{\mathbb{E} [\|\mathbf{X}_t - \mathbf{x}^*\|]} = \exp\left(-\frac{c}{n}\right)$$

$$\mathbb{E} \log \frac{\|\mathbf{X}_{t+1} - \mathbf{x}^*\|}{\|\mathbf{X}_t - \mathbf{x}^*\|} = -\frac{c}{n}$$

Connection with Hitting Time formulation

$$T_\epsilon \approx \frac{n}{c} \log \frac{\epsilon_0}{\epsilon}$$

Convergence Rates - Hitting time - Wrap up

	Rate of convergence	Hitting time scaling
Pure Random Search (1+1)-ES constant step-size	$\frac{1}{t} \log \frac{\ \mathbf{X}_t - \mathbf{x}^*\ }{\ \mathbf{X}_0 - \mathbf{x}^*\ } \approx -\frac{1}{n} \frac{\log(t)}{t}$	$\left(\frac{\epsilon_0}{\epsilon}\right)^n$
Linear Convergence (fixed n) + Linear dependence wrt n	$\mathbb{E} [\ \mathbf{X}_t - \mathbf{x}^*\] = \exp\left(-\frac{c}{n}\right)^t \mathbb{E} [\ \mathbf{X}_0 - \mathbf{x}^*\]$ $\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\ \mathbf{X}_t - \mathbf{x}^*\ }{\ \mathbf{X}_0 - \mathbf{x}^*\ } = -\frac{c}{n}$	$\frac{n}{c} \log \frac{\epsilon_0}{\epsilon}$