

# Covariance Matrix Adaptation Evolution Strategy (CMA-ES)

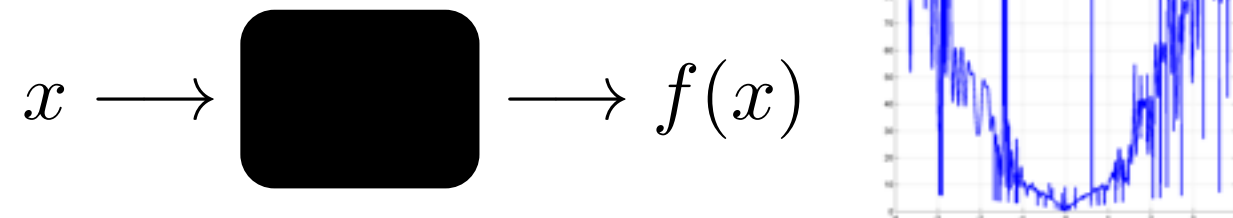
Nikolaus Hansen  
Inria  
Research Centre Saclay  
Machine Learning and Optimization Team (TAO)  
Univ. Paris-Sud, LRI

# Problem Statement: Black-Box Optimization

Given an objective function

$$f : \mathcal{X} \subset \mathbb{R}^n \rightarrow \mathbb{R}$$

Minimize  $f$  in a *black-box scenario* (direct search, no gradients)



Problem domain specific knowledge is only used *within* the black-box

## Objective

- **convergence** to a global essential infimum of  $f$  as fast as possible  
linear convergence,  $\mathcal{O}(n \log 1/\epsilon)$  black-box evaluations
- **find**  $x \in \mathcal{X}$  with small  $f(x)$  value using as few black-box calls as possible

The black box can

- be non-convex, multi-modal/rugged, discontinuous, noisy, dynamic
- take from milli-seconds to hours to evaluate

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*



# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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* ↻

# Landscape of Continuous Search Methods

## *Gradient-based (Taylor, local)*

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

## *Derivative-free optimization (DFO)*

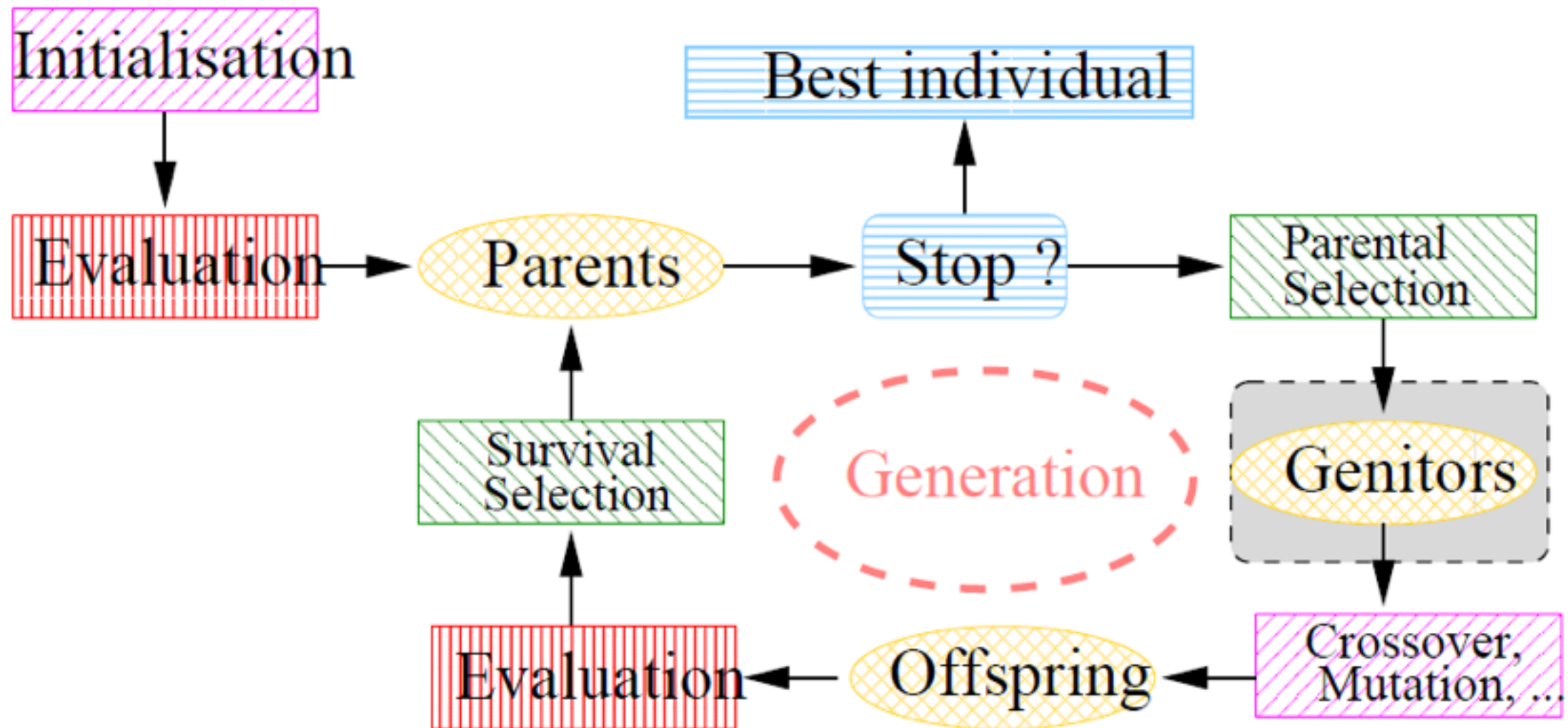
- **Trust-region methods** (NEWUOA, BOBYQA) [Powell 2006, 2009]
- **Simplex downhill** [Nelder & Mead 1965]
- **Pattern search** [Hooke & Jeeves 1961, Audet & Dennis 2006]





## *Stochastic (randomized) search methods*

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



# A Different Viewpoint: Evolutionary Algorithm Scheme



-  Stochastic operators Representation dependent
-  Darwinian Evolution Engine (can be stochastic or deterministic)
-  Main CPU cost
-  Checkpointing: stopping criterion, statistics, updates, ...

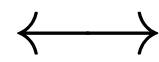


# Metaphors

## Evolutionary Computation

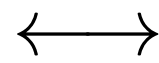
## Optimization/Nonlinear Programming

individual, offspring, parent



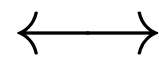
candidate solution  
 decision variables  
 design variables  
 object variables

population



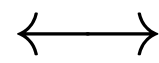
set of candidate solutions

fitness function



objective function  
 loss function  
 cost function  
 error function

generation



iteration

...methods: ESs

# The CMA-ES

**Input:**  $\mathbf{m} \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ ,  $\lambda$

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

**Set:**  $c_c \approx 4/n$ ,  $c_\sigma \approx 4/n$ ,  $c_1 \approx 2/n^2$ ,  $c_\mu \approx \mu_w/n^2$ ,  $c_1 + c_\mu \leq 1$ ,  $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$ ,  
and  $w_{i=1\dots\lambda}$  such that  $\mu_w = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \approx 0.3 \lambda$

**While not terminate**

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

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

$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbb{1}_{\{\|\mathbf{p}_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w$  cumulation for  $\mathbf{C}$

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$  update  $\mathbf{C}$

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

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

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Evolution Strategies

New search points are sampled normally distributed

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

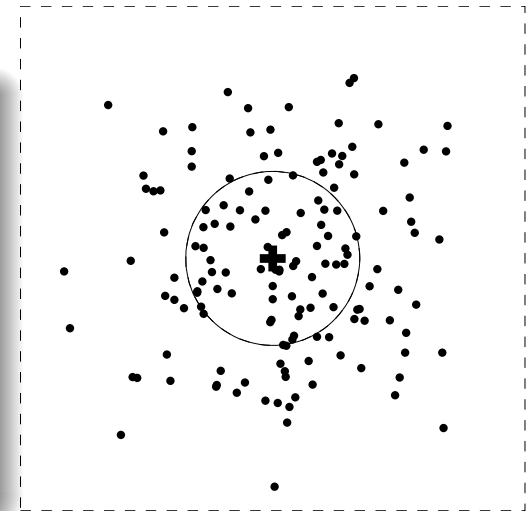
as perturbations of  $\mathbf{m}$ , where  $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ ,  $\mathbf{C} \in \mathbb{R}^{n \times n}$

where

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

here, all new points are sampled with the same parameters

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



# Evolution Strategies

New search points are sampled normally distributed

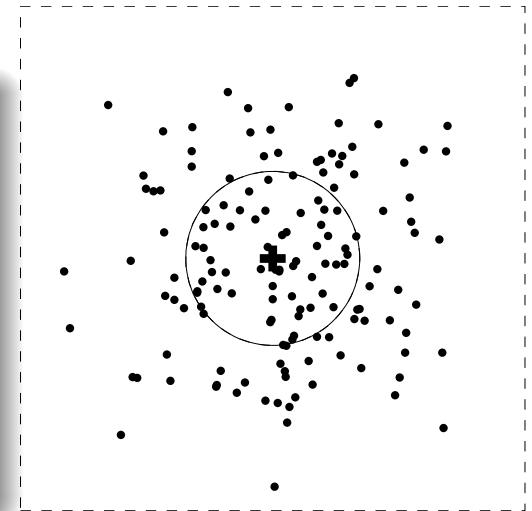
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

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

where

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

here, all new points are sampled with the same parameters



The (multivariate) normal distribution (Gaussian distribution)

# 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:

$$\mathcal{N}(\mathbf{x}, \mathbf{A}) + \mathcal{N}(\mathbf{y}, \mathbf{B}) \sim \mathcal{N}(\mathbf{x} + \mathbf{y}, \mathbf{A} + \mathbf{B})$$

helpful in **design and analysis** of algorithms  
related to the *central limit theorem*

3 most convenient way to generate **isotropic** search points

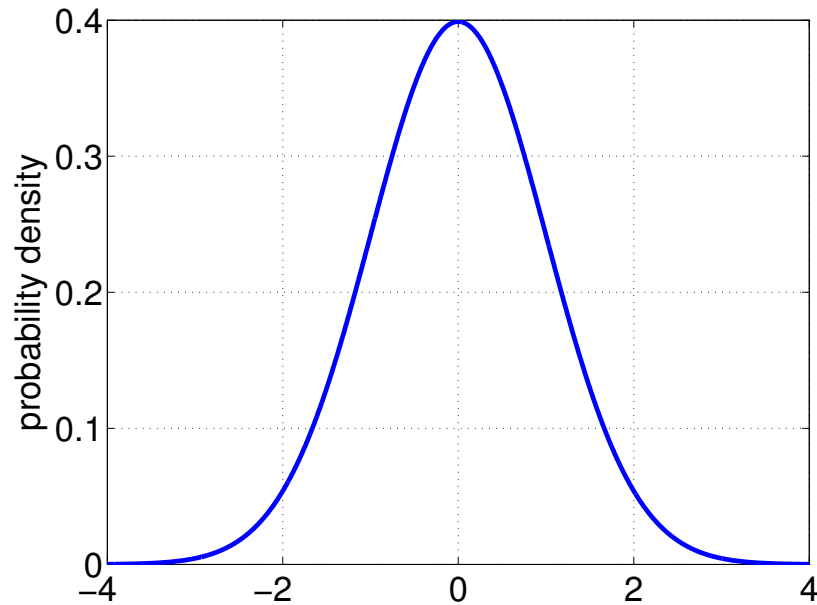
the isotropic distribution does **not favor any direction**, rotational invariant

4 maximum entropy distribution with finite variance

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

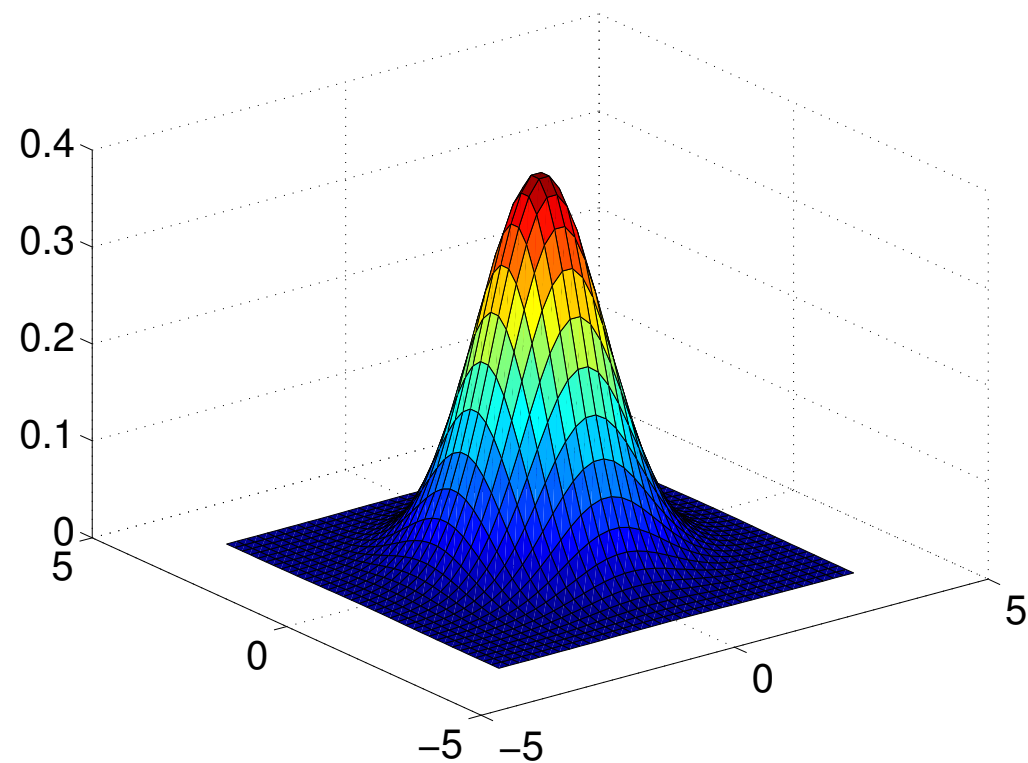
# Normal Distribution

Standard Normal Distribution

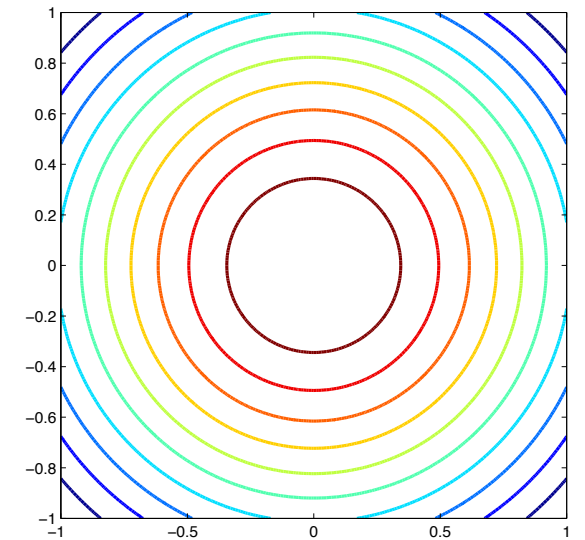


probability density of the 1-D standard normal distribution

2-D Normal Distribution



probability density of a 2-D normal distribution



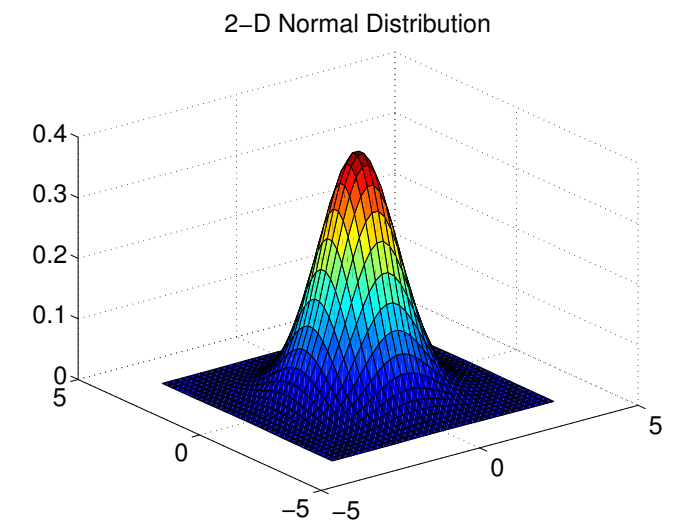


# The Multi-Variate ( $n$ -Dimensional) Normal Distribution

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

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

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

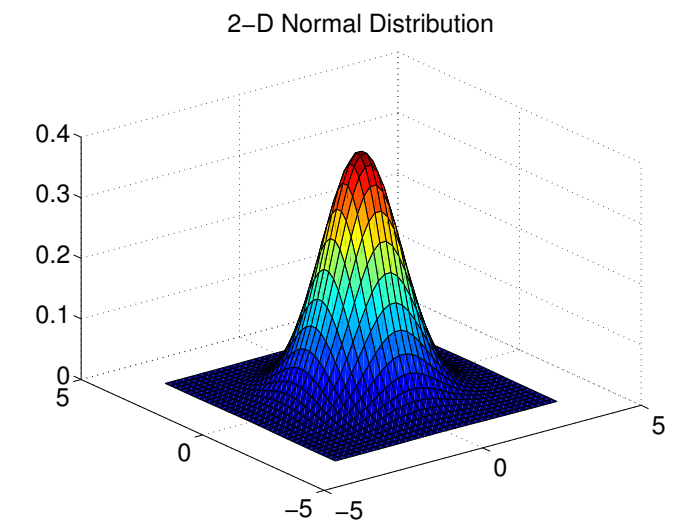


# The Multi-Variate ( $n$ -Dimensional) Normal Distribution

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

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

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

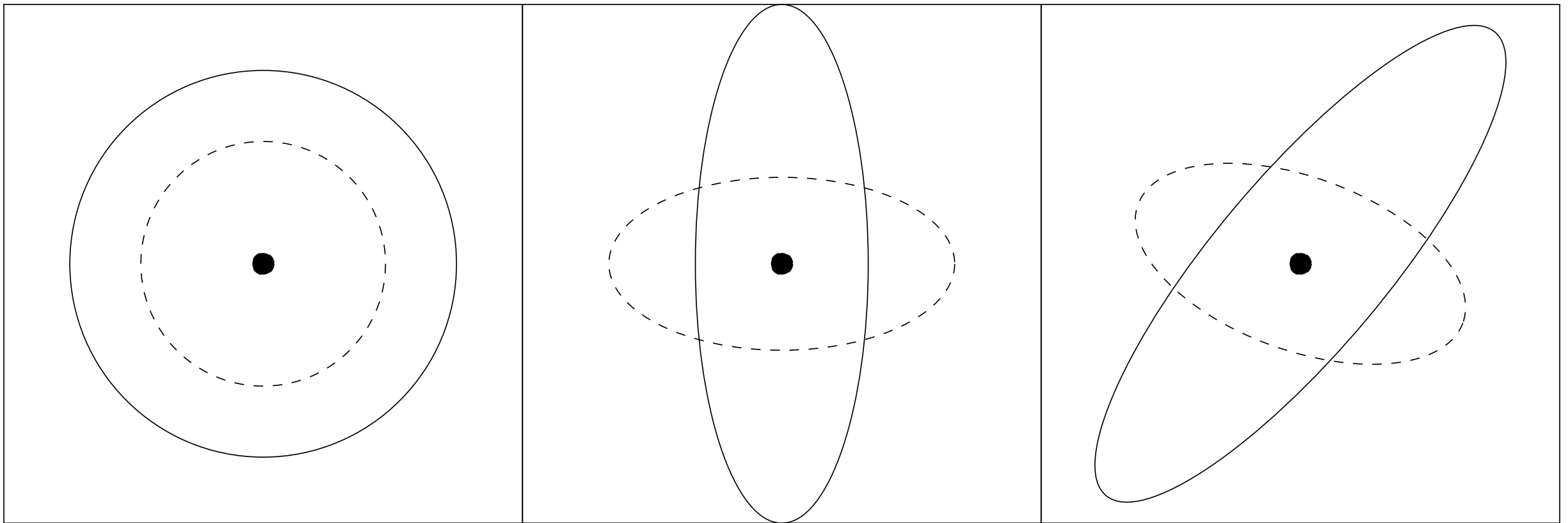


The **covariance matrix**  $\mathbf{C}$

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

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

### Lines of Equal Density



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

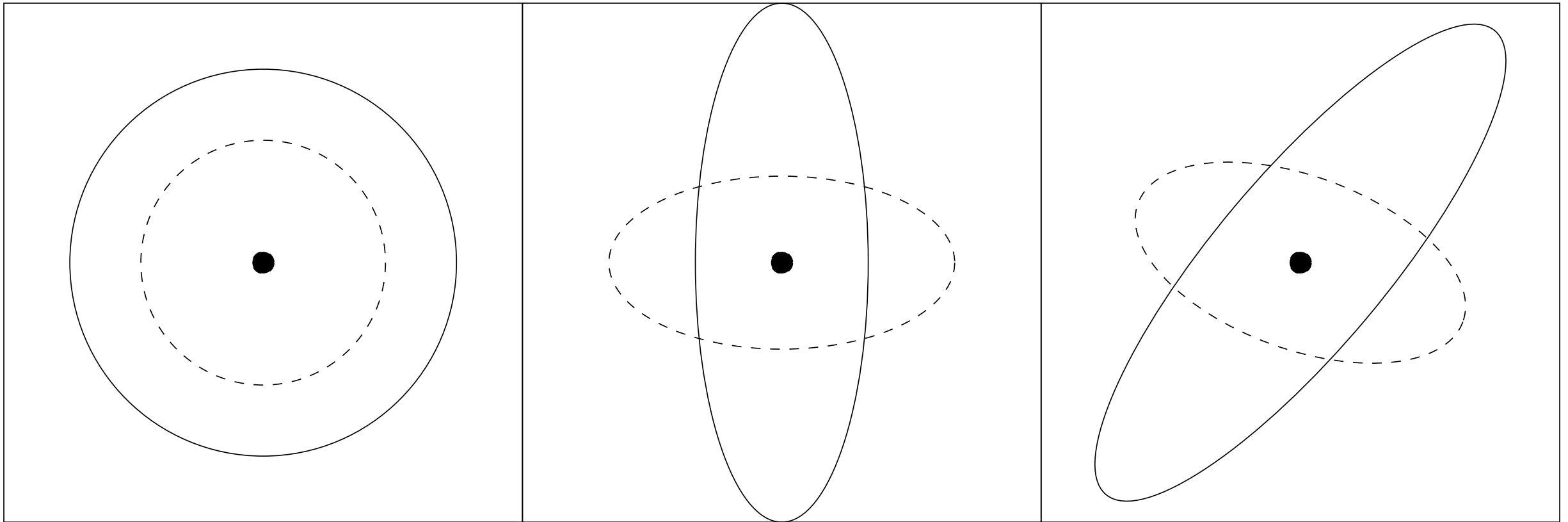
$\mathcal{N}(\mathbf{m}, \mathbf{D}^2) \sim \mathbf{m} + \mathbf{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$   
 **$n$  degrees of freedom**  
 components are  
 independent, scaled

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

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

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

### Lines of Equal Density



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

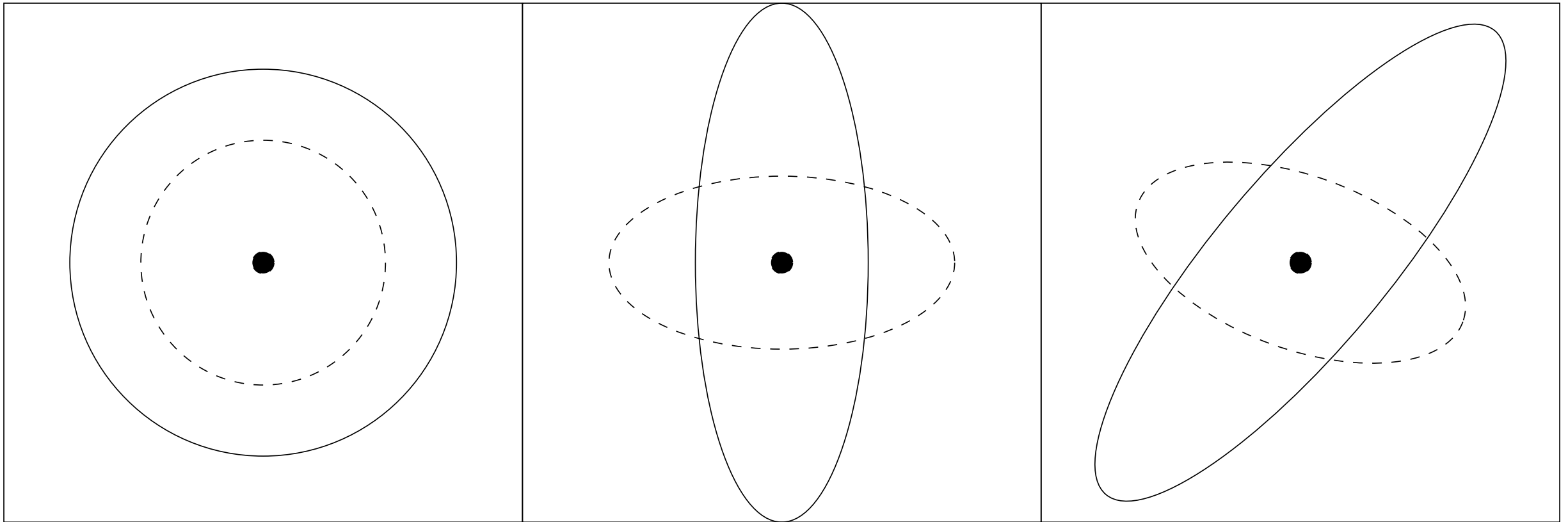
$\mathcal{N}(\mathbf{m}, \mathbf{D}^2) \sim \mathbf{m} + \mathbf{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$   
 **$n$  degrees of freedom**  
 components are  
 independent, scaled

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

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

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

### Lines of Equal Density



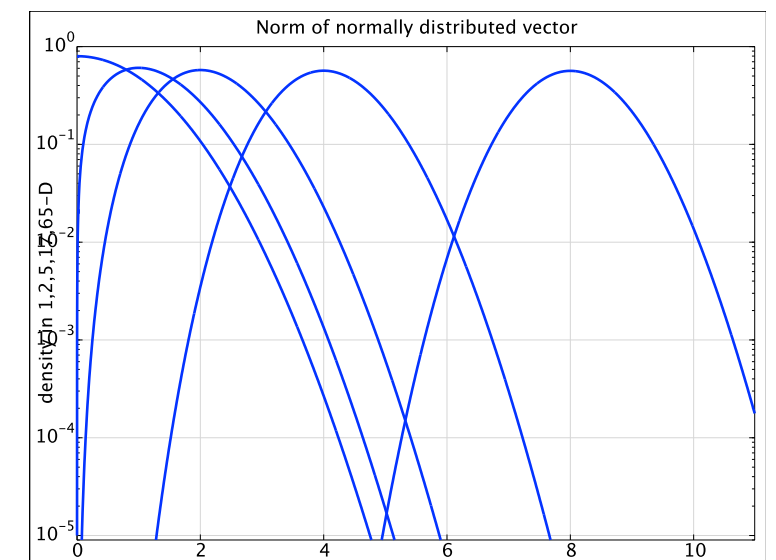
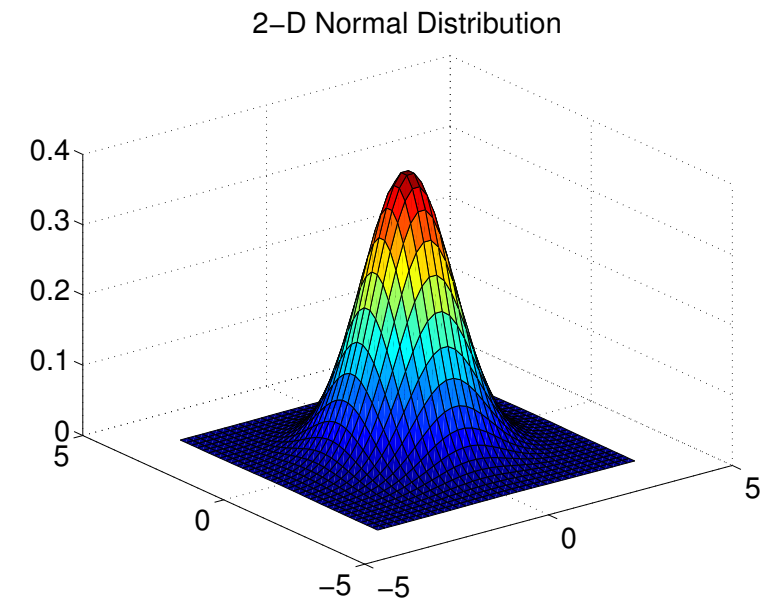
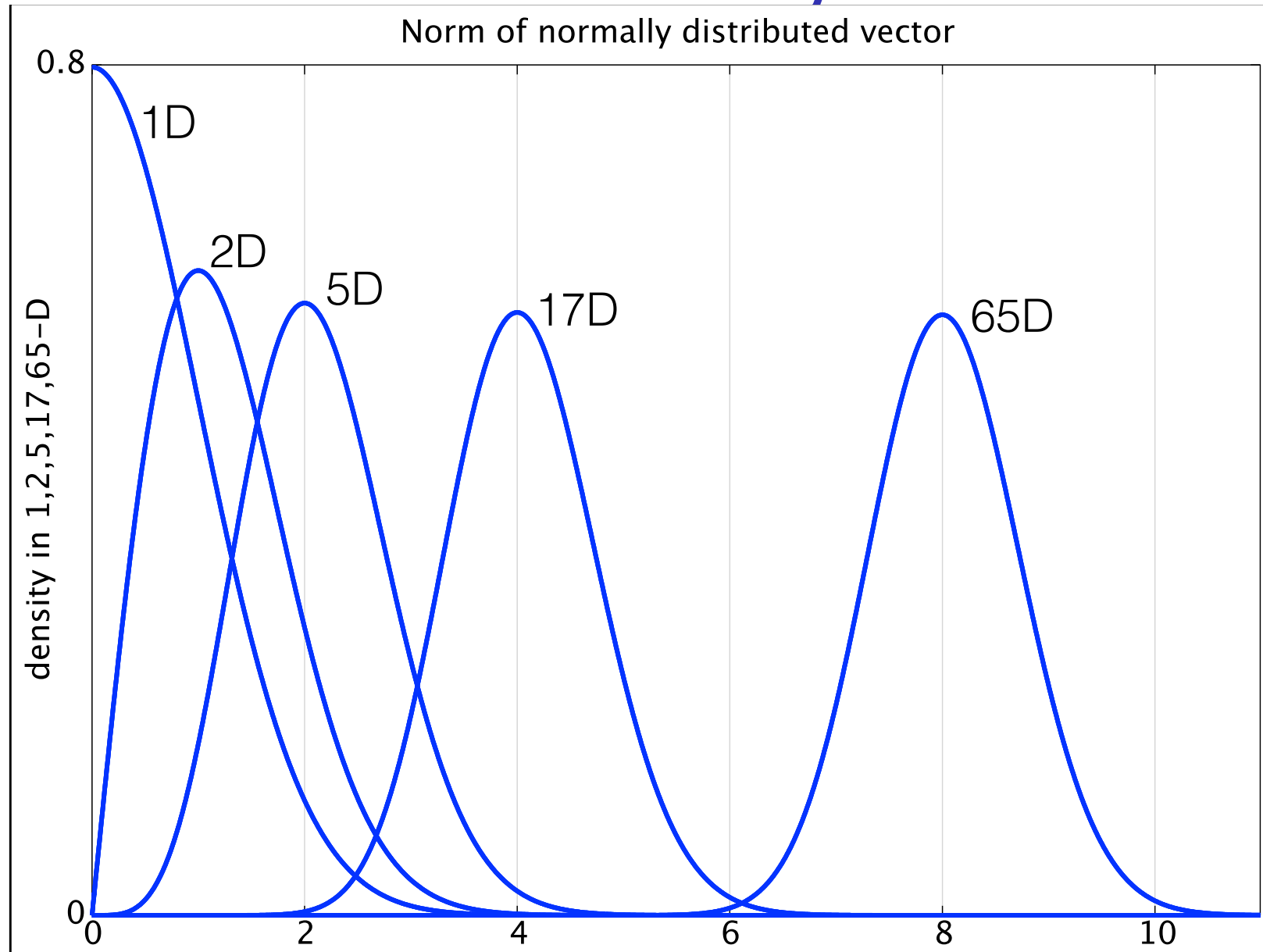
$\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{I}) \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{I})$   
**one degree of freedom**  $\sigma$   
 components are  
 independent standard  
 normally distributed

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

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

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

# Effect of Dimensionality



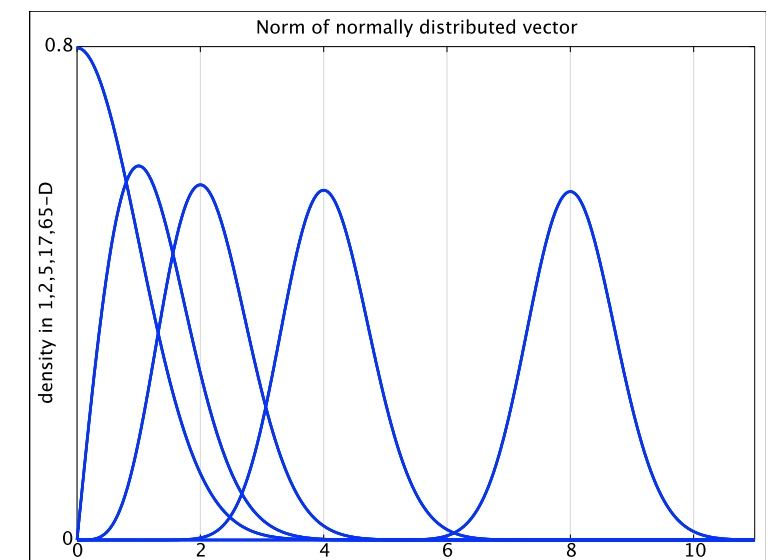
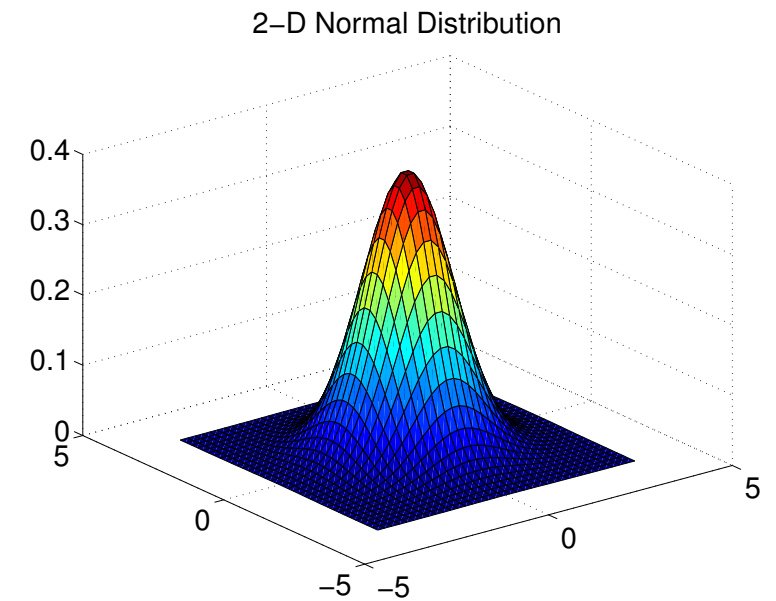
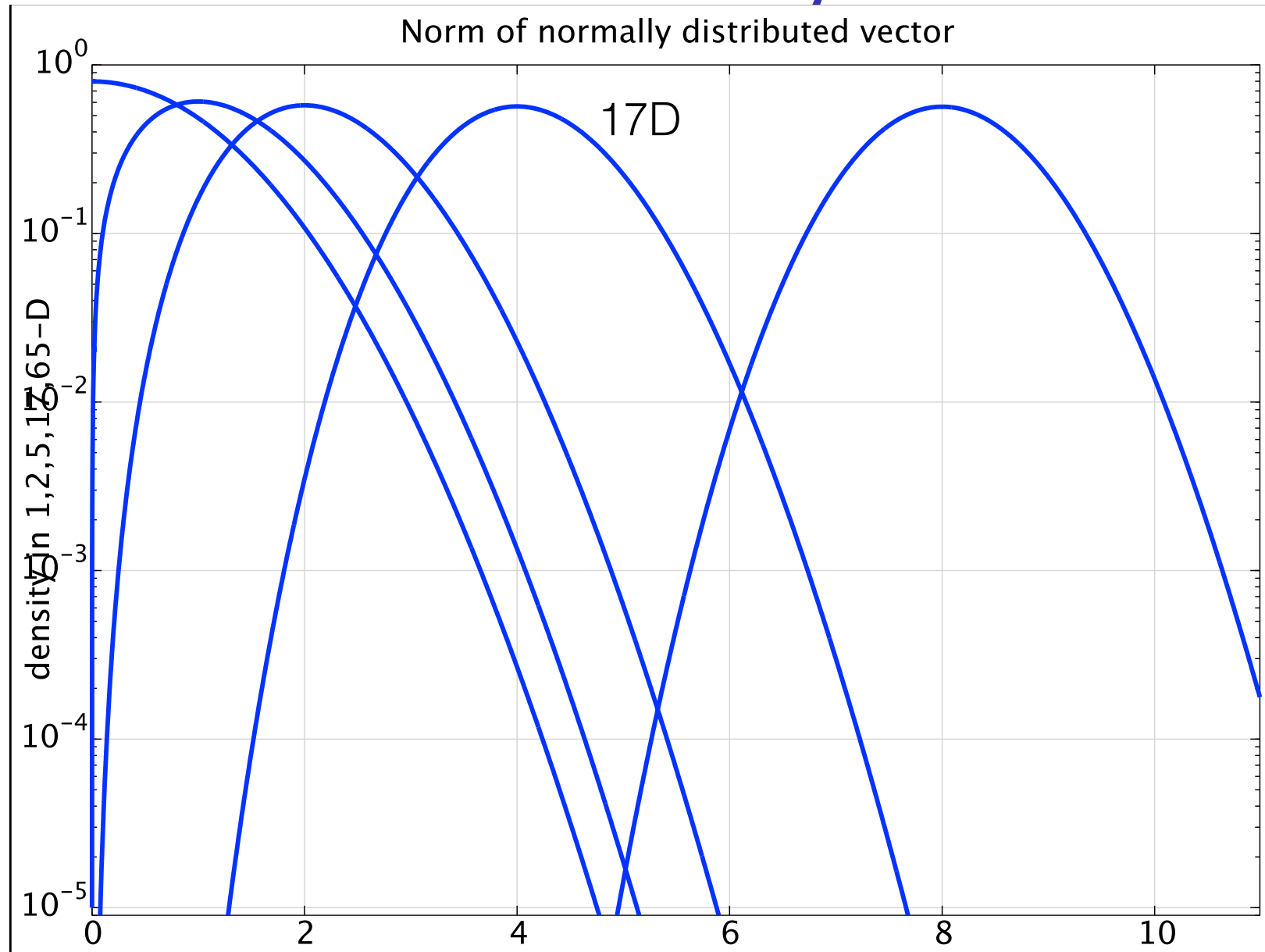
$\|\mathcal{N}(\mathbf{0}, \mathbf{I})\| \longrightarrow \mathcal{N}\left(\sqrt{n-1/2}, 1/2\right)$  with modal value  $\sqrt{n-1}$

yet: maximum entropy distribution

also consider a difference between two vectors:

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

# Effect of Dimensionality



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

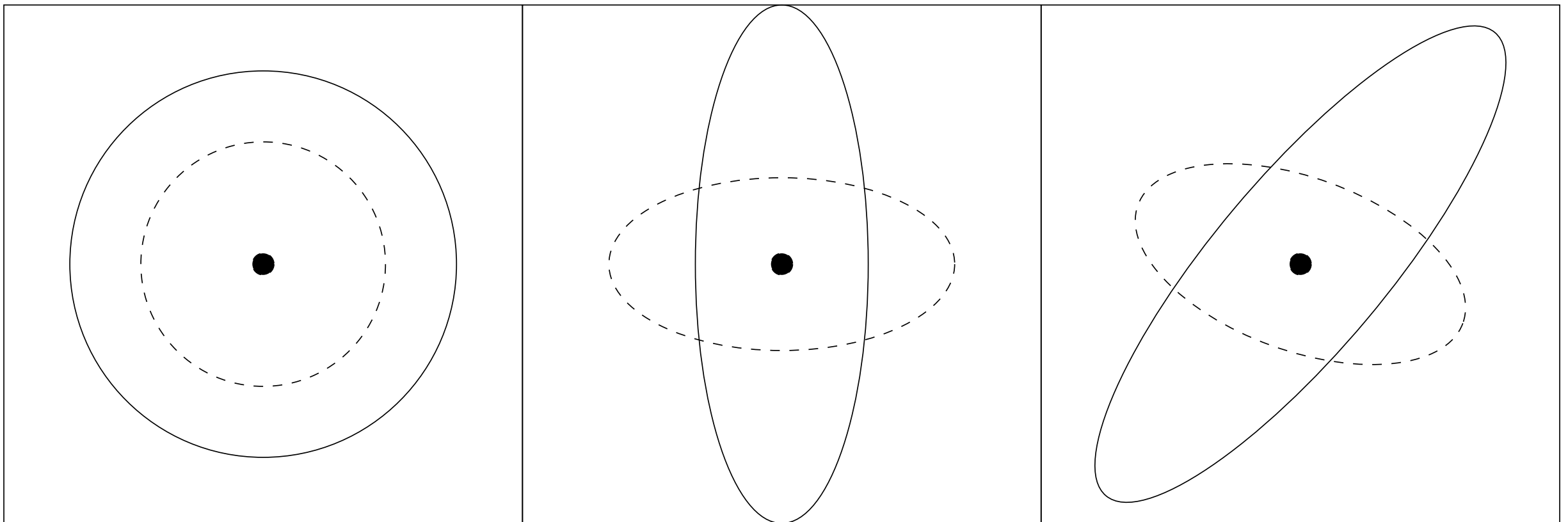
yet: maximum entropy distribution

also consider a difference between two vectors:

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

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

Lines of Equal Density



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



Update of the distribution mean

# Stochastic Search

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

Initialize distribution parameters  $\theta$ , 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  $F_\theta$

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*

# Evolution Strategies

New search points are sampled normally distributed

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

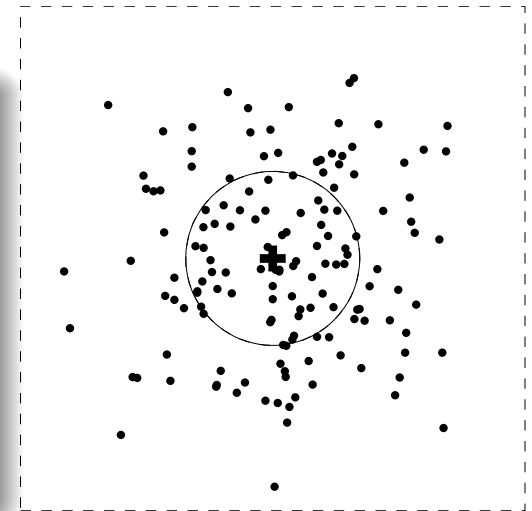
as perturbations of  $\mathbf{m}$ , where  $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ ,  $\mathbf{C} \in \mathbb{R}^{n \times n}$

where

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

here, all new points are sampled with the same parameters

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



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

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$$

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

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

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$$

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

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

Non-elitist selection and intermediate (weighted) recombination

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

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

The new mean reads

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

where

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

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

# Evolution Strategies

## Terminology

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

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

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

$(\mu, \lambda)$ -ES: 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

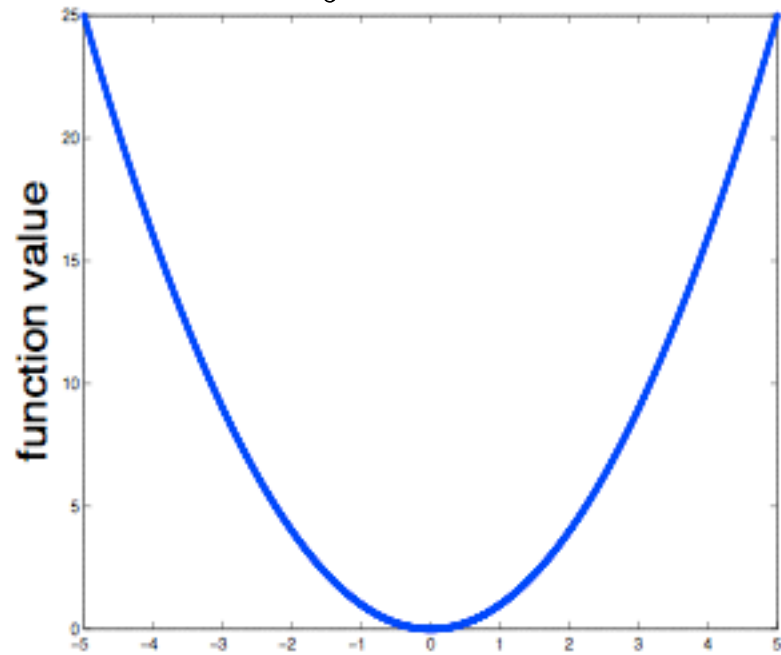
$$m \leftarrow x$$

Invariance

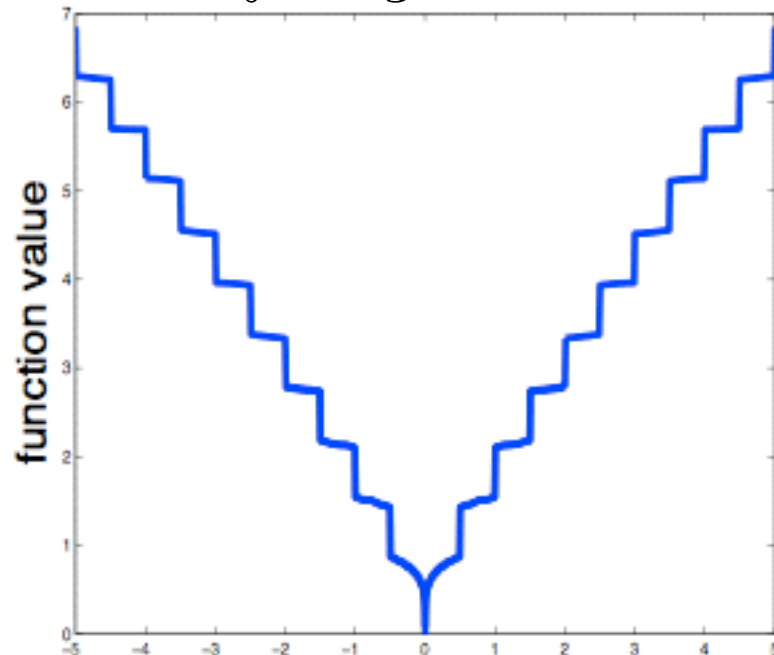


# Invariance: Function-Value Free Property

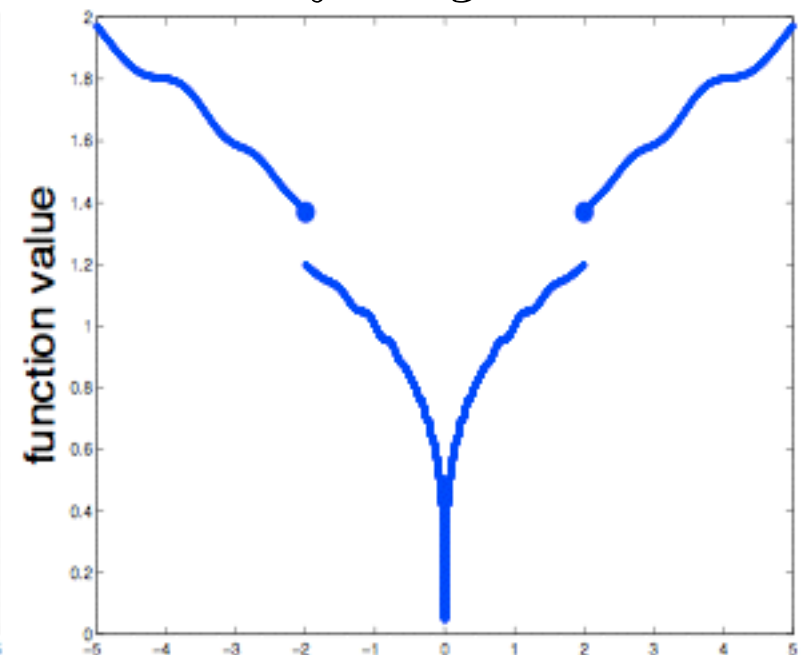
$$f = h$$



$$f = g_1 \circ h$$



$$f = g_2 \circ h$$



Three functions belonging to the same equivalence class

A *function-value free search algorithm* is invariant under the transformation with any **order preserving** (strictly increasing)  $g$ .

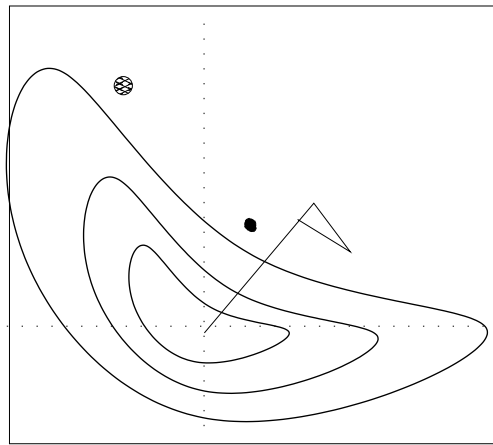
Invariances make

- observations meaningful as a rigorous notion of generalization
- algorithms predictable and/or "robust"

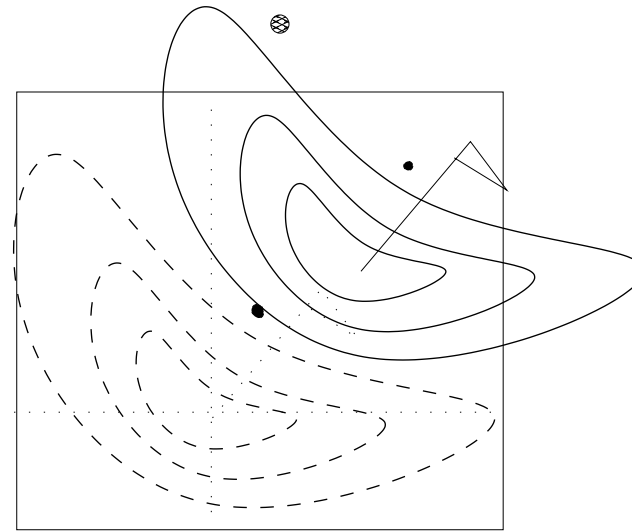
# Basic Invariance in Search Space

- translation invariance

is true for most optimization algorithms



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



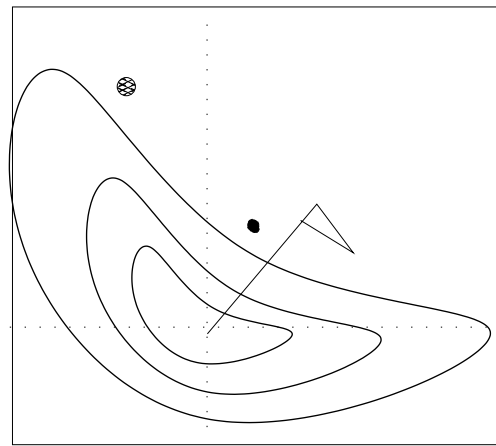
Identical behavior on  $f$  and  $f_a$

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

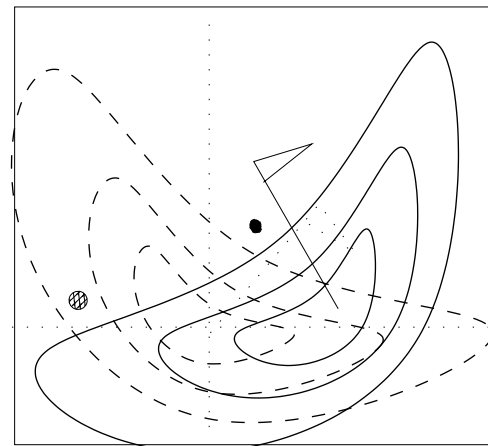
No difference can be observed w.r.t. the argument of  $f$

# Rotational Invariance in Search Space

- invariance to orthogonal (rigid) transformations  $\mathbf{R}$ , where  $\mathbf{R}\mathbf{R}^T = \mathbf{I}$   
 e.g. true for simple evolution strategies  
 recombination operators might jeopardize rotational invariance



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



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

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

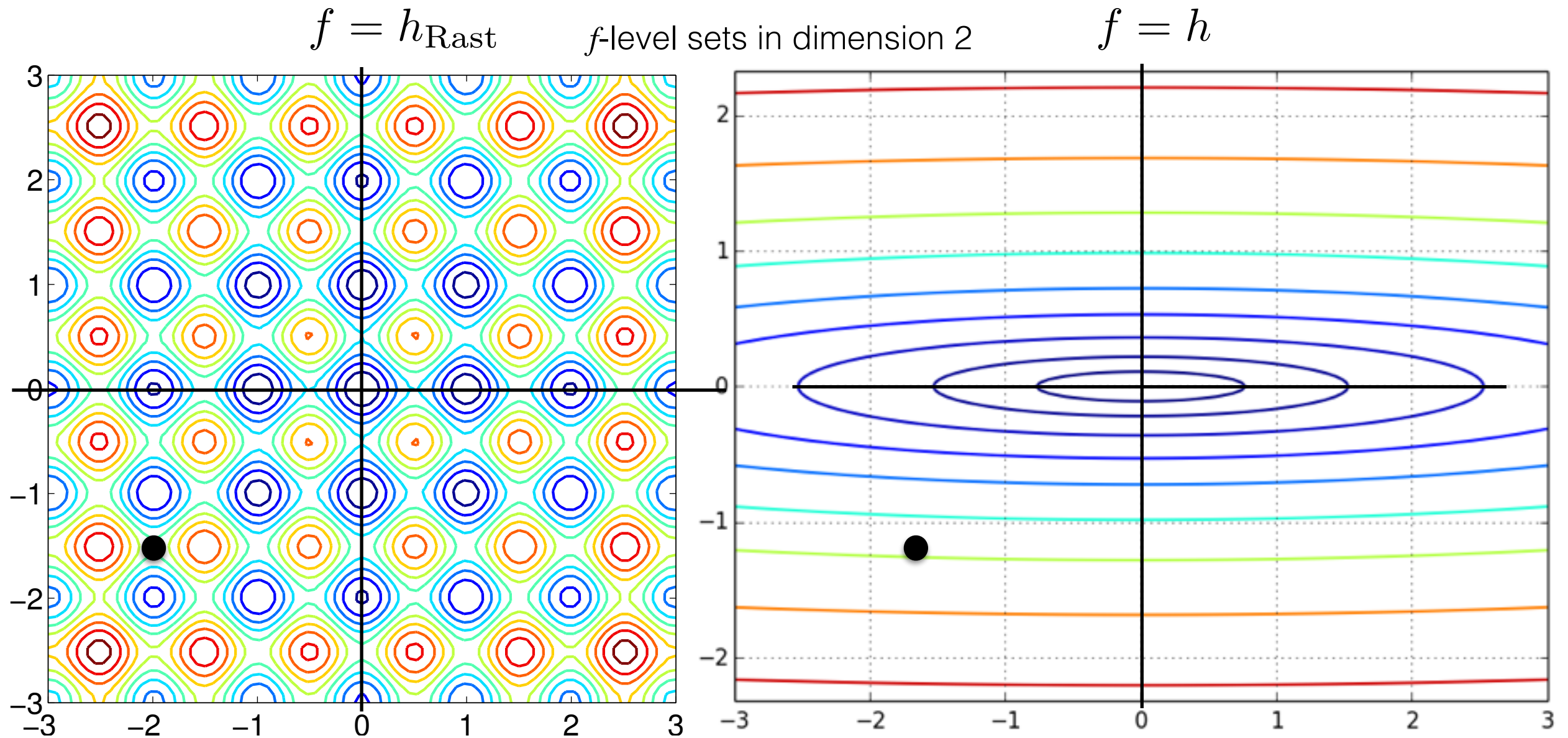
45

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

<sup>4</sup> 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

<sup>5</sup> Hansen 2000. Invariance, Self-Adaptation and Correlated Mutations in Evolution Strategies. *Parallel Problem Solving from Nature PPSN VI*

# Invariance Under Rigid Search Space Transformations



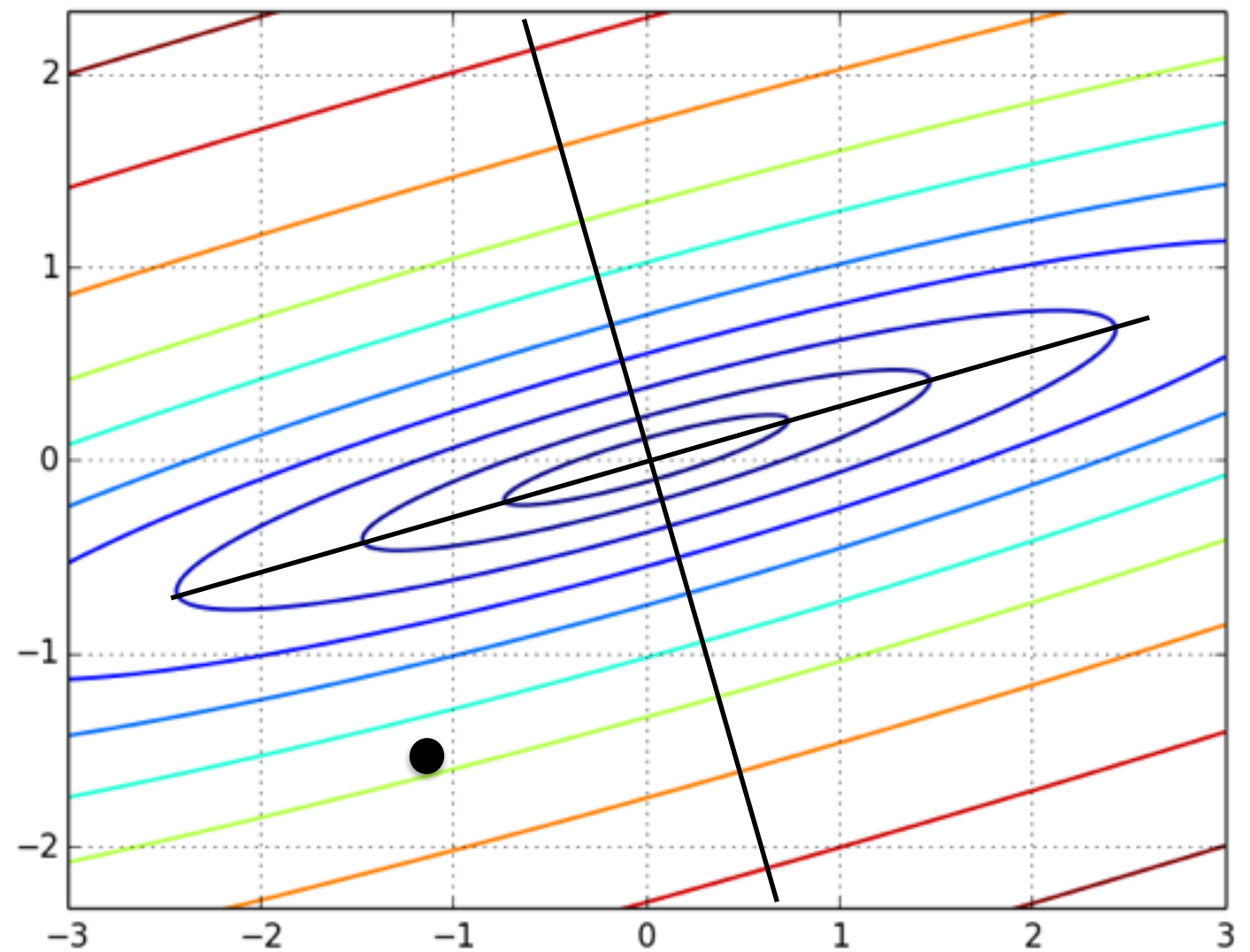
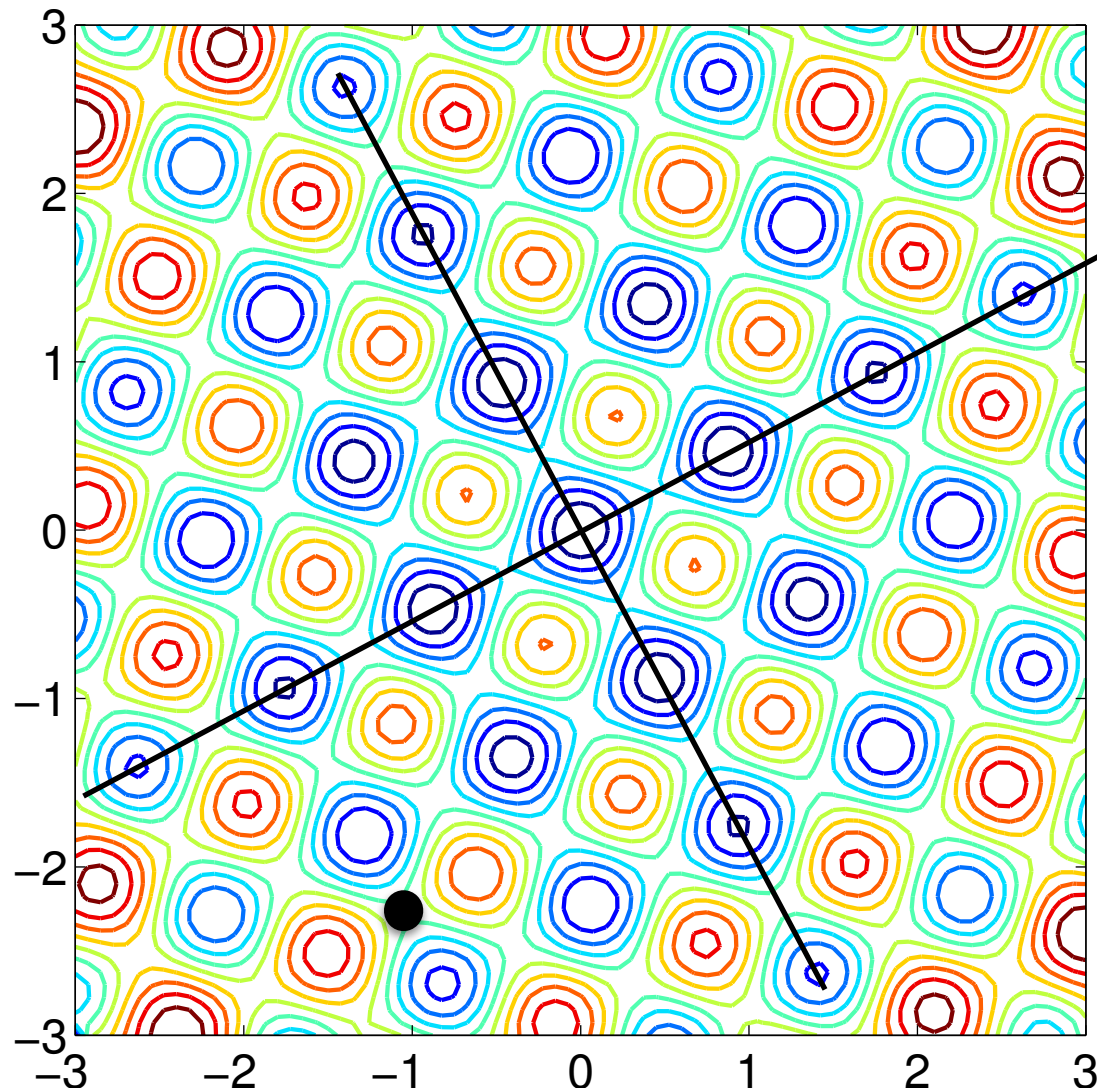
for example, invariance under search space rotation  
 (separable  $\Leftrightarrow$  non-separable)



# Invariance Under Rigid Search Space Transformations

$$f = h_{\text{Rast}} \circ R \quad f\text{-level sets in dimension 2}$$

$$f = h \circ R$$



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

# Landscape of Continuous Search Methods

## *Gradient-based (Taylor, local)*

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

## *Derivative-free optimization (DFO)*

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

## *Stochastic (randomized) search methods*

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

# Invariance

*The grand aim of all science is to cover the greatest number of empirical facts by logical deduction from the smallest number of hypotheses or axioms.*

— Albert Einstein

- Empirical performance results
  - ▶ from benchmark functions
  - ▶ from solved real world problems

are only useful if they do **generalize** to other problems

- **Invariance** is a strong **non-empirical** statement about generalization
  - generalizing (identical) performance from a single function to a whole class of functions

consequently, invariance is important for the evaluation of search algorithms

# Step-Size Control



# Evolution Strategies

Recalling

New search points are sampled normally distributed

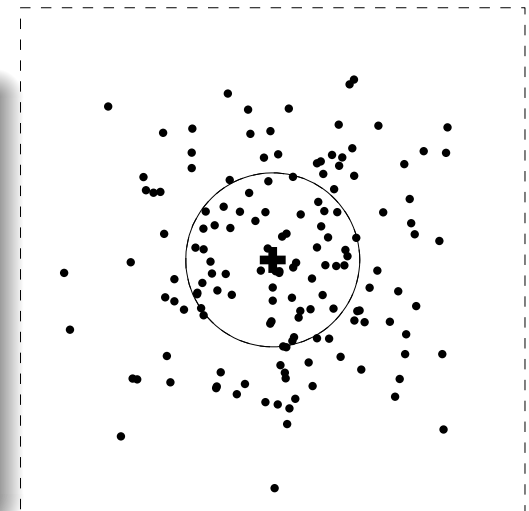
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

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

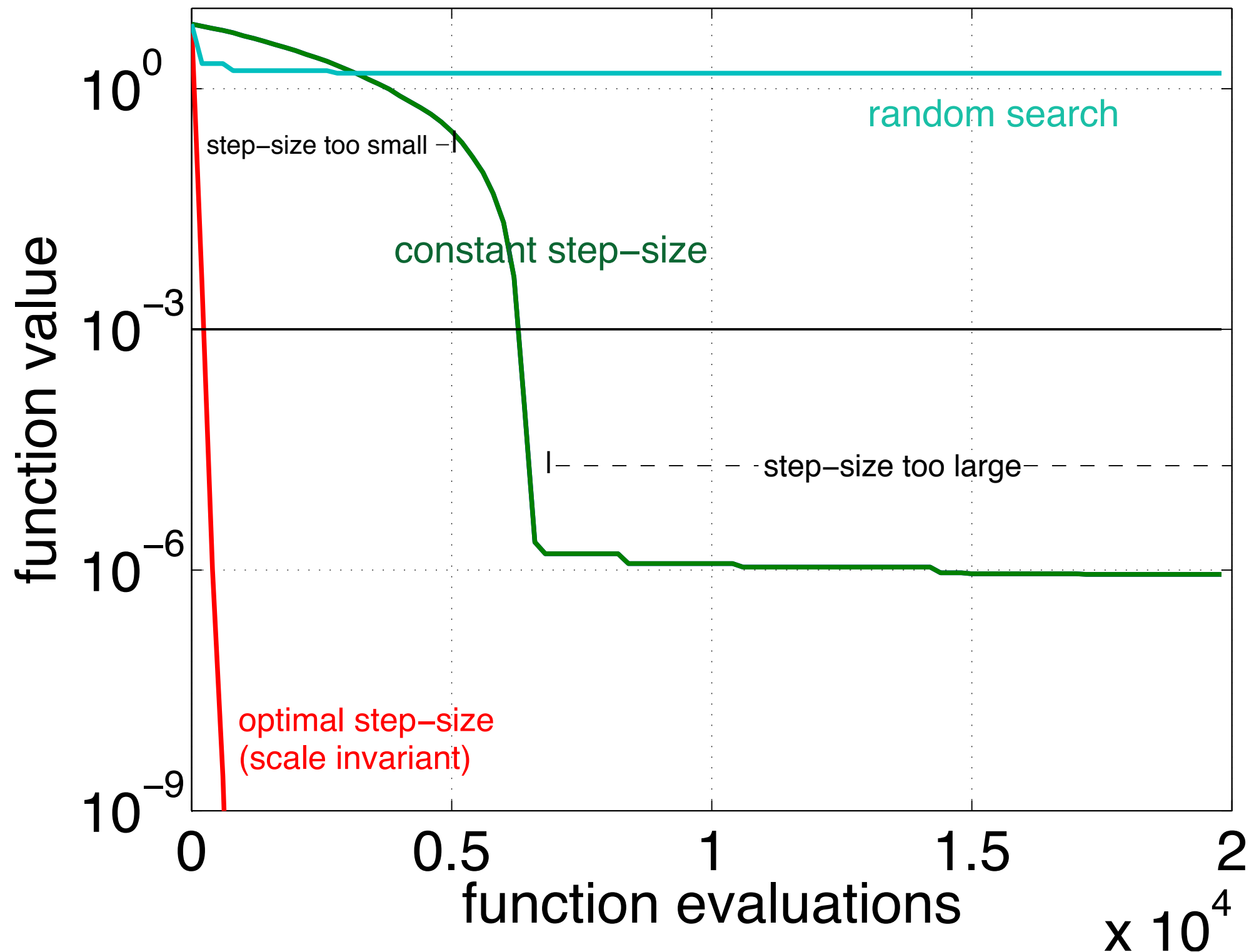
where

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

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



# Why Step-Size Control?



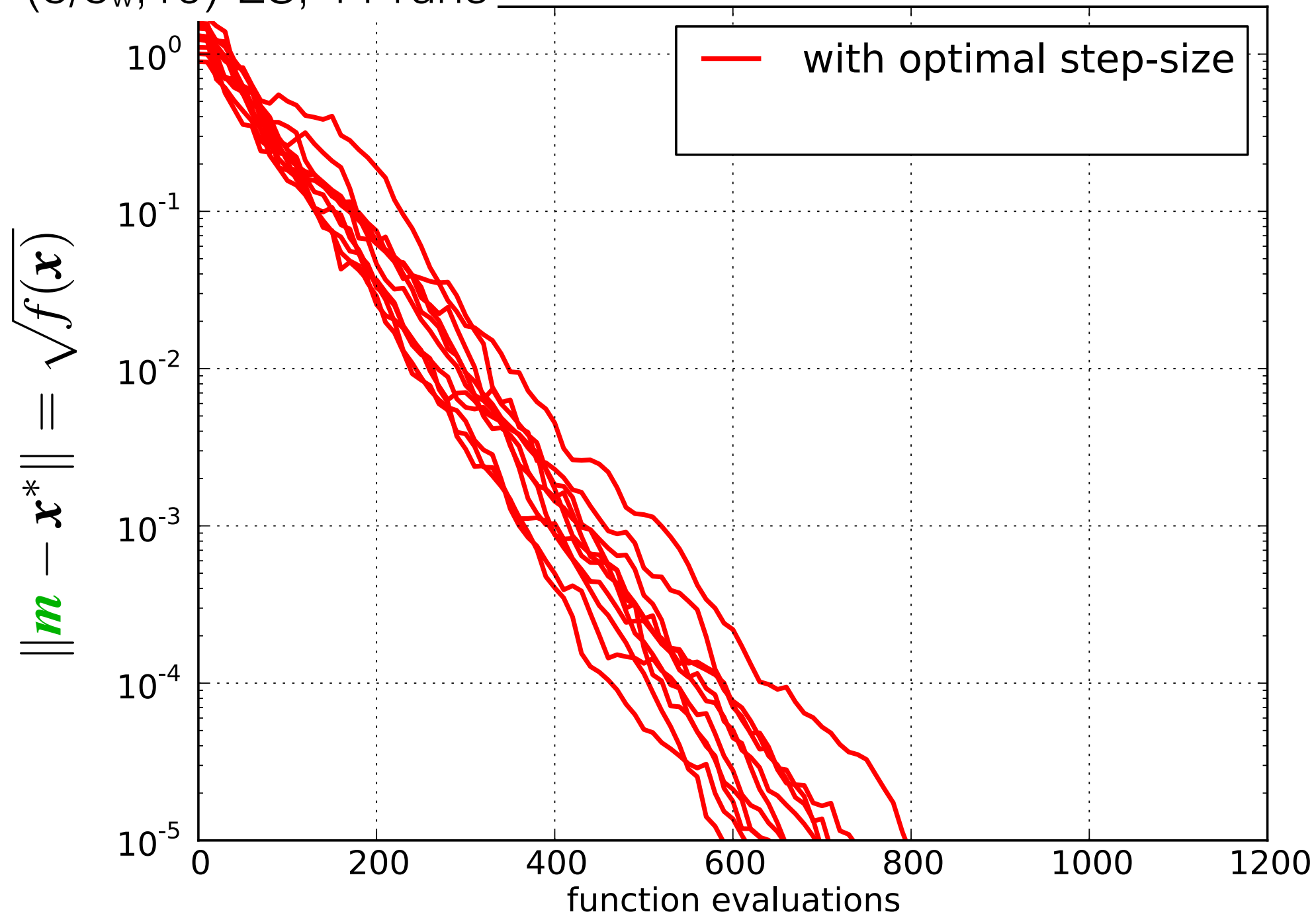
(1+1)-ES  
(red & green)

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

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

# Why Step-Size Control?

(5/5<sub>w</sub>,10)-ES, 11 runs



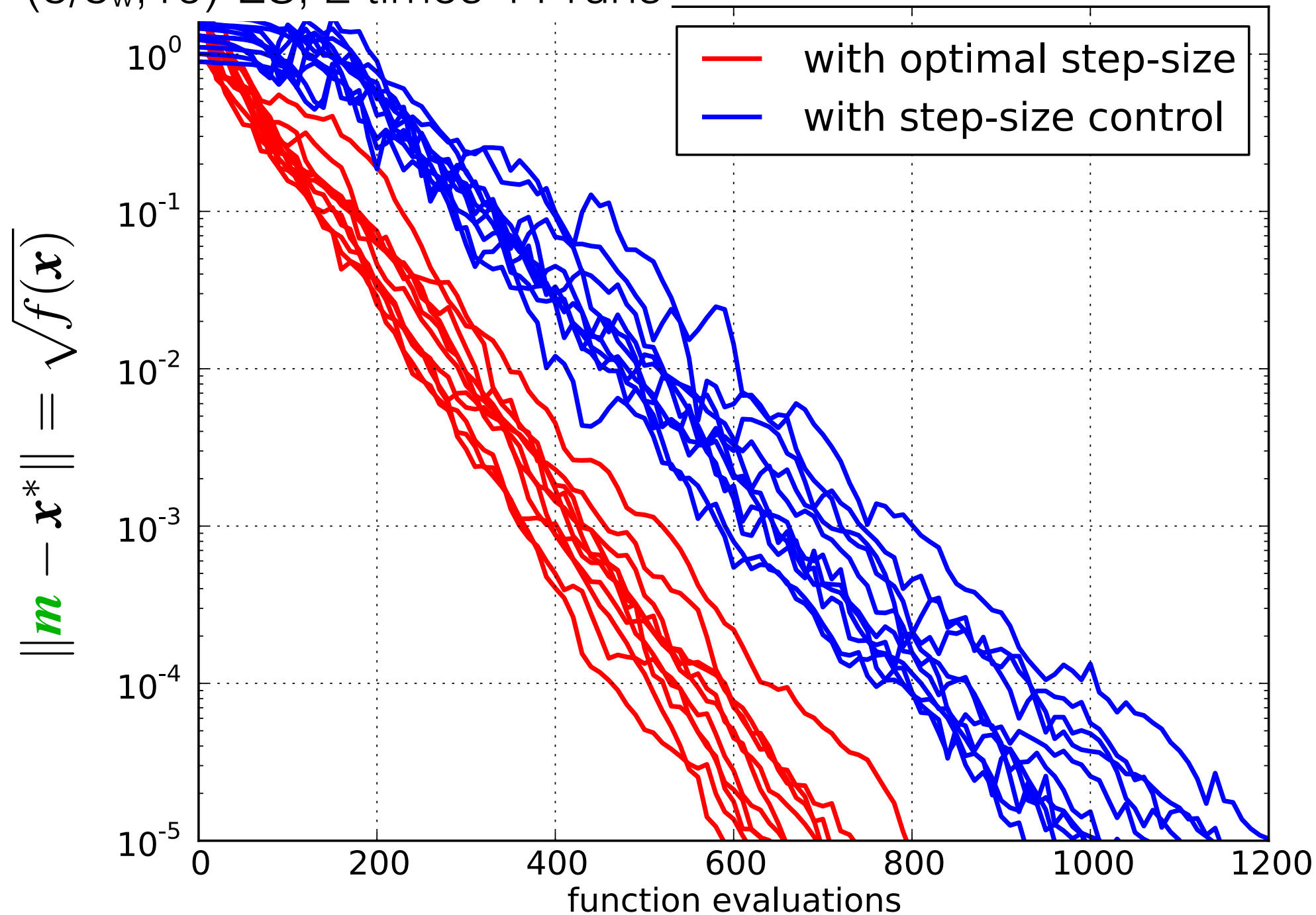
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

with optimal step-size  $\sigma$

# Why Step-Size Control?

(5/5<sub>w</sub>,10)-ES, 2 times 11 runs



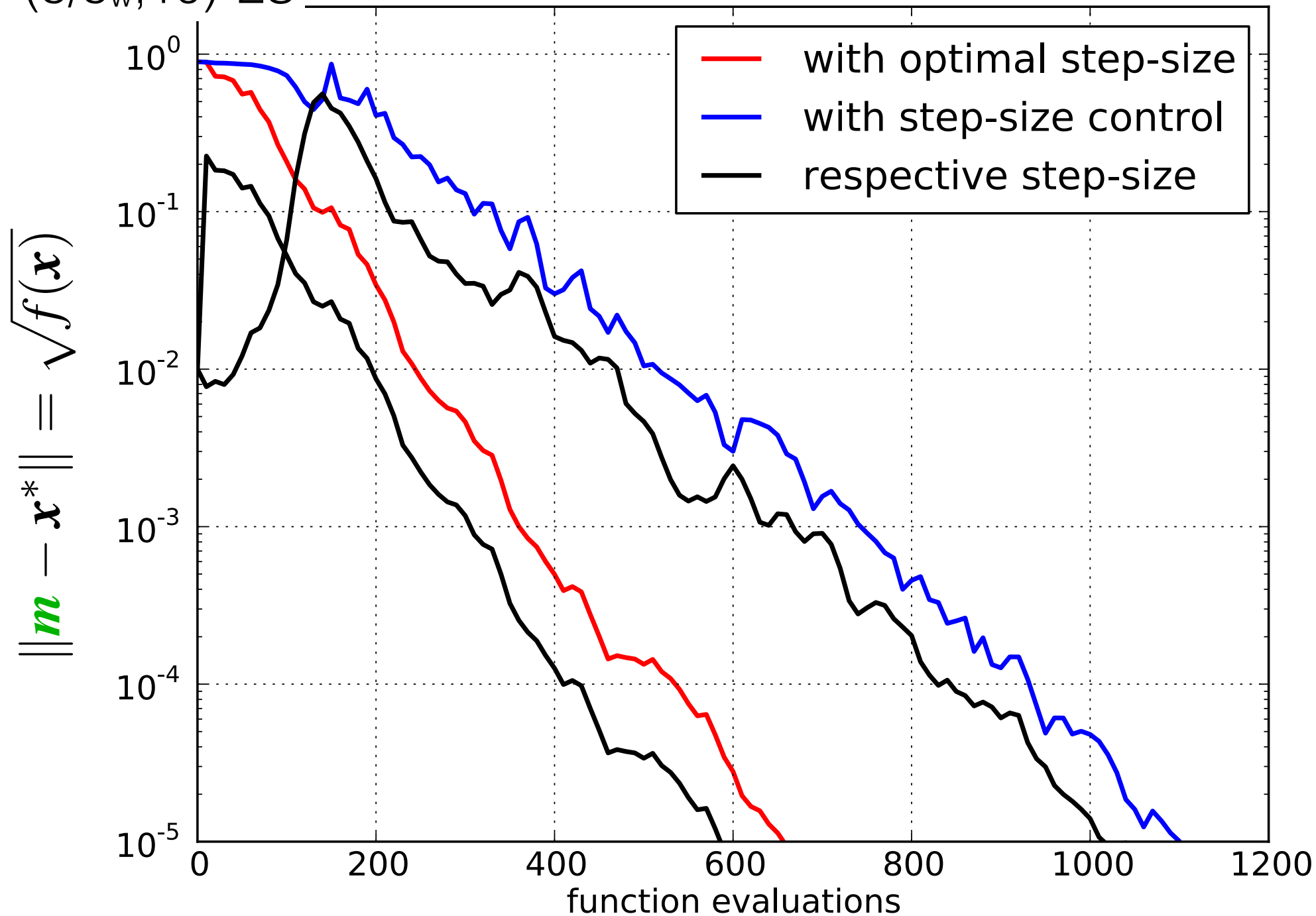
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

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

# Why Step-Size Control?

(5/5<sub>w</sub>, 10)-ES



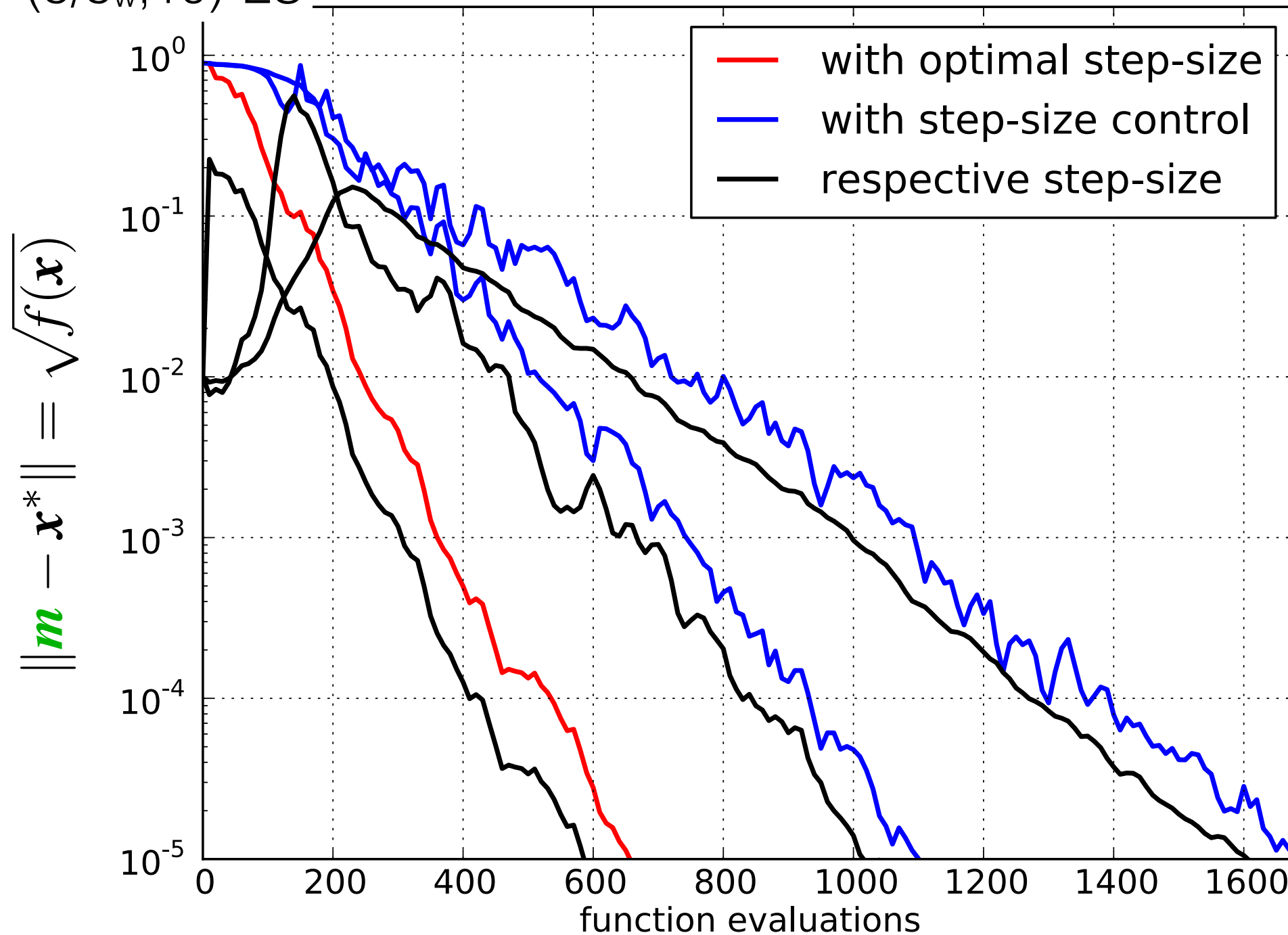
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

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

# Why Step-Size Control?

(5/5<sub>w</sub>, 10)-ES



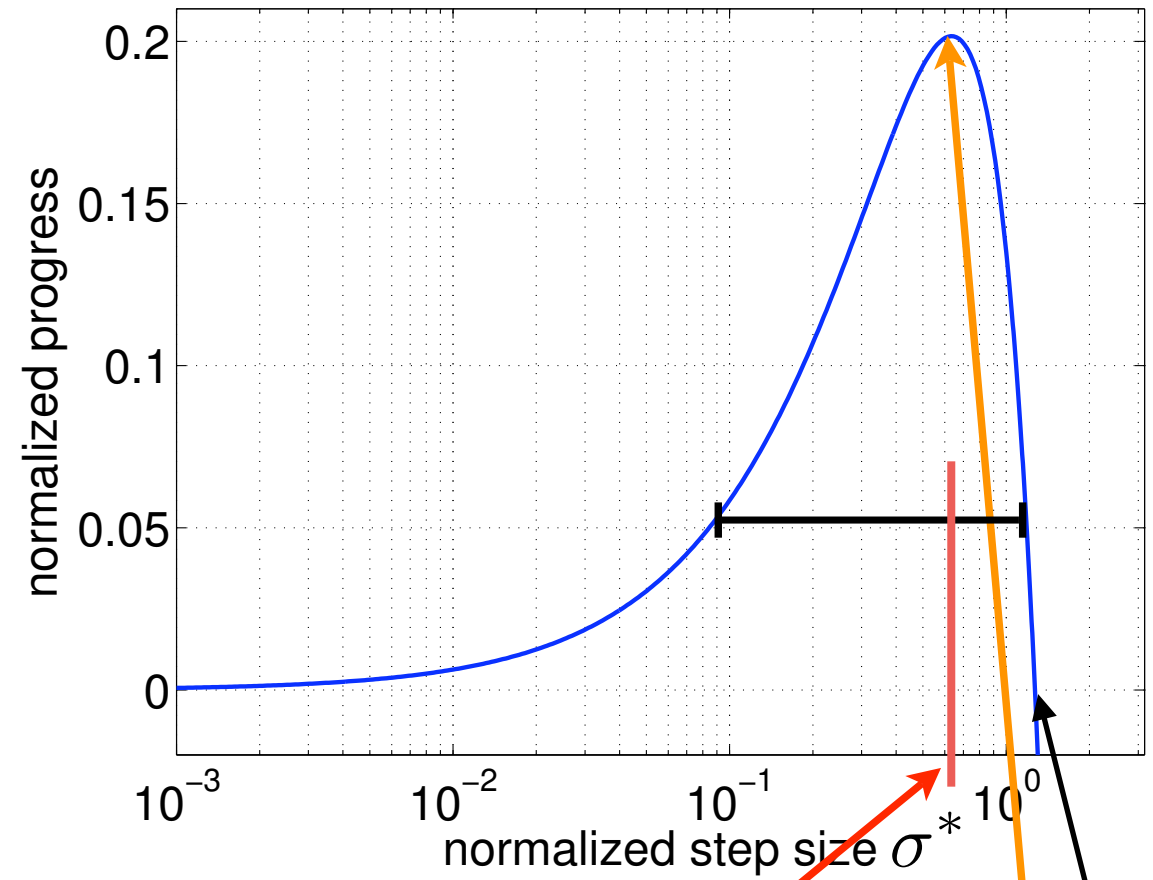
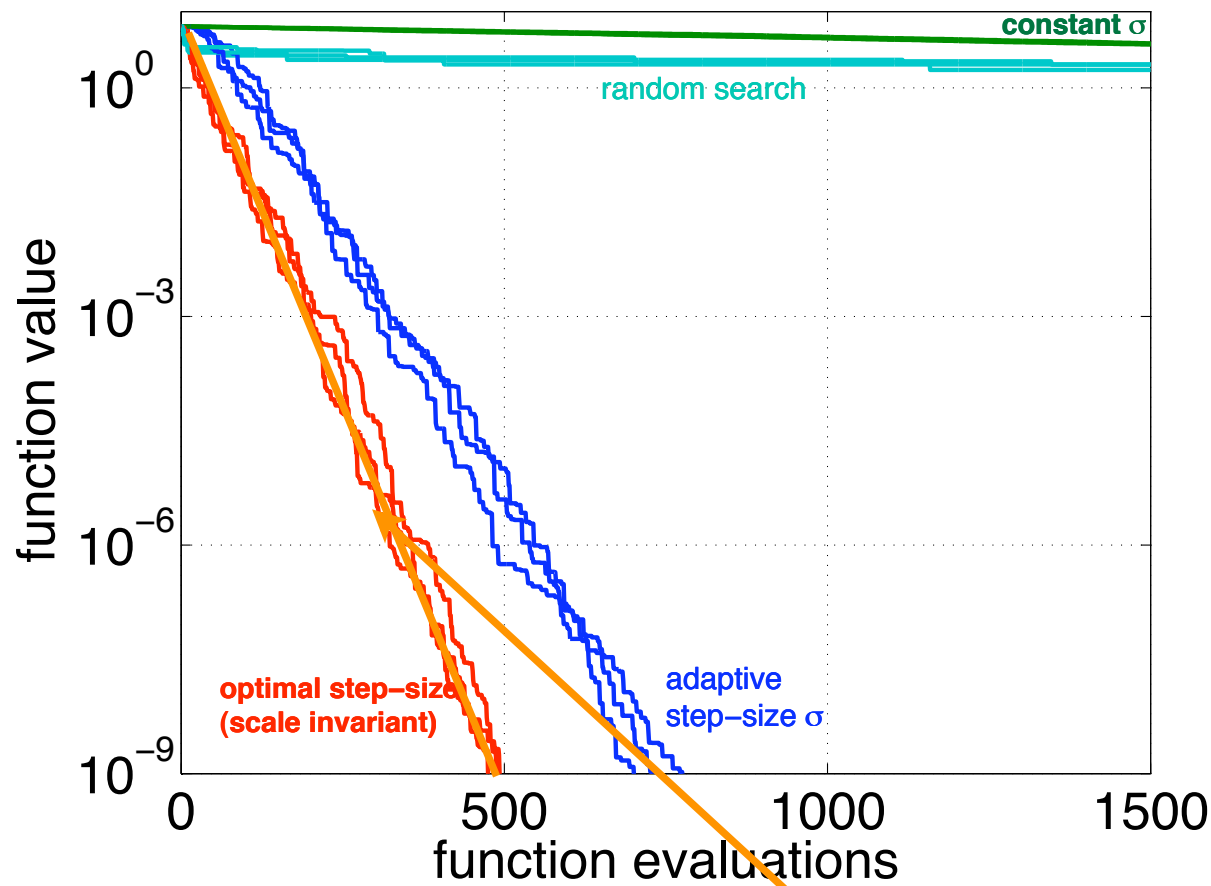
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

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

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

# Why Step-Size Control?

$$\sigma_{\text{opt}} \approx \sigma_{\text{opt}}^* \mu_w \frac{\|m\|}{n}$$



$$f(m^{(t)}) \approx f(m^{(0)}) \times \exp\left(2 \frac{-\varphi^*}{n} \lambda t\right)$$

$$\frac{\lg(f(m^{(t+1)})) - \lg(f(m^{(t)}))}{\lambda} \approx \lg(e) \times 2 \frac{-\varphi^*}{n}$$

$\sigma_{\text{opt}}^*$   $\varphi^*$   
 transition from convergence to divergence

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

# Methods for Step-Size Control

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

---

<sup>a</sup>Rechenberg 1973, *Evolutionsstrategie, Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*, Frommann-Holzboog

<sup>b</sup>Schumer and Steiglitz 1968. Adaptive step size random search. *IEEE TAC*

<sup>c</sup>Schwefel 1981, *Numerical Optimization of Computer Models*, Wiley

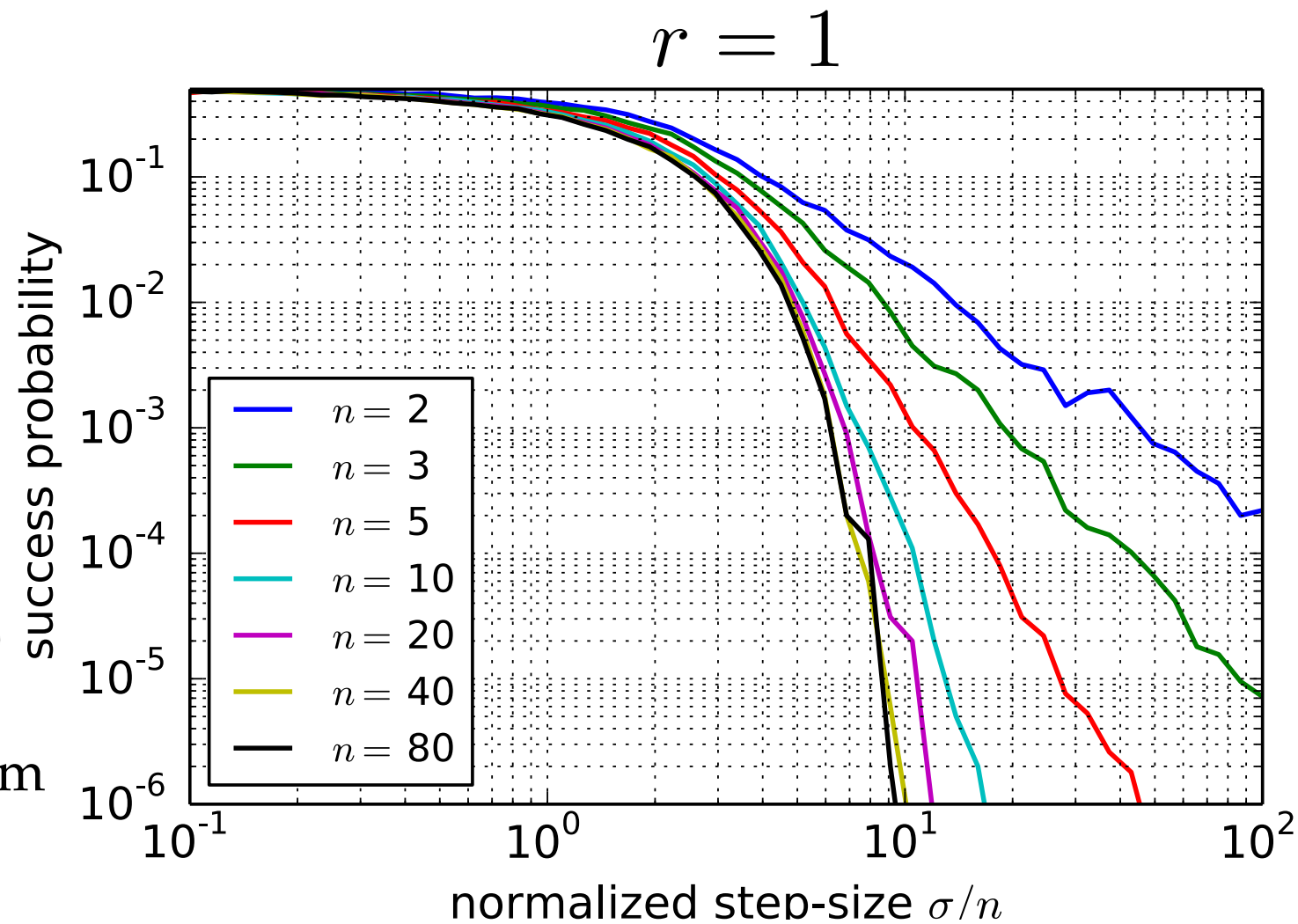
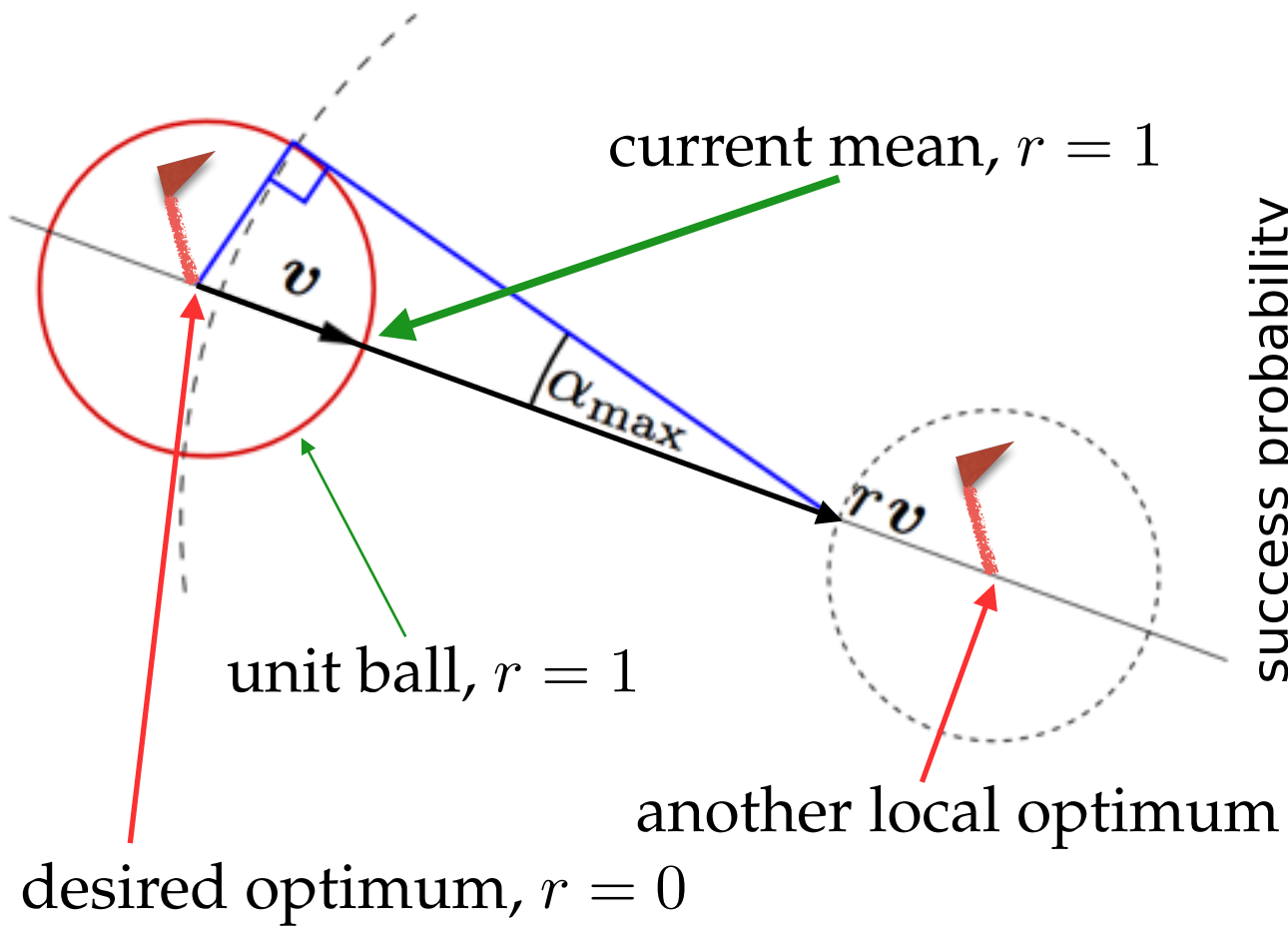
<sup>d</sup>Hansen & Ostermeier 2001, Completely Derandomized Self-Adaptation in Evolution Strategies, *Evol. Comput.*

9(2)

<sup>e</sup>Ostermeier *et al* 1994, Step-size adaptation based on non-local use of selection information, *PPSN IV*

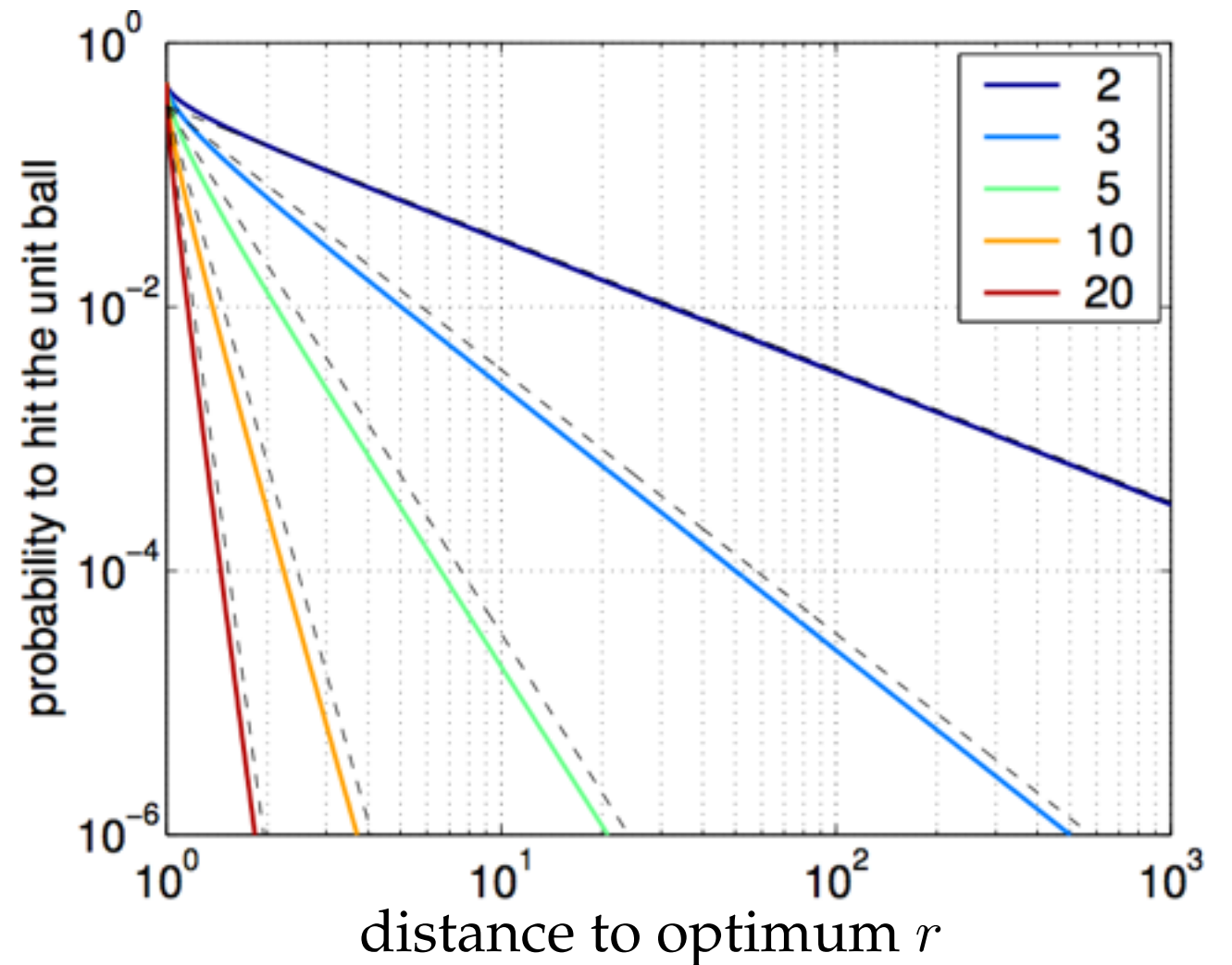
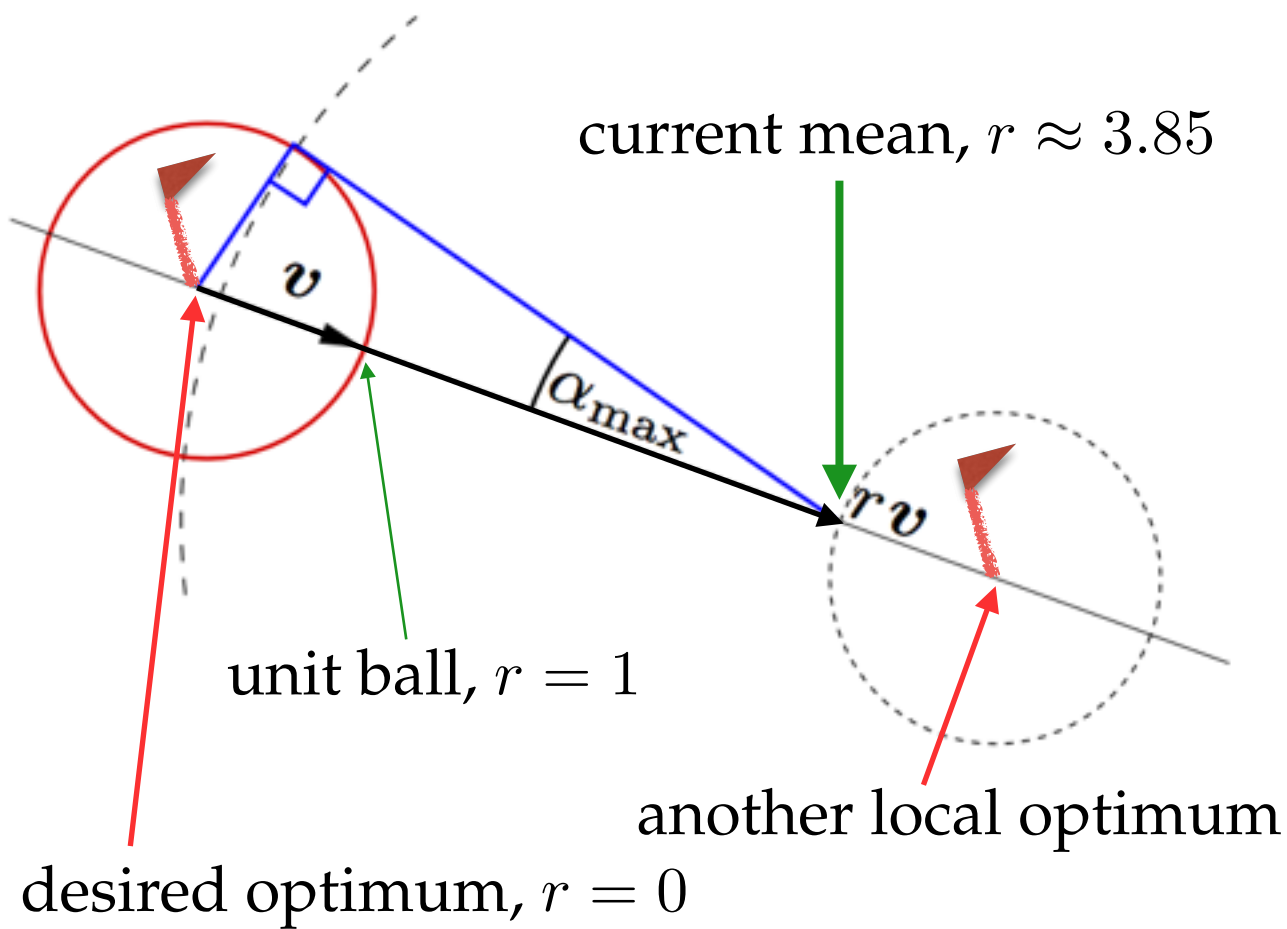


# Success vs step-size



probability to hit the unit ball as a function of step-size

# Success vs distance



**Fig. 1.** Probability to hit the unit hyperball (solid) sampling from  $r\mathbf{v}$  as mean with an optimal isotropic distribution, where  $\mathbf{v} \in \mathbb{R}^n$  and  $\|\mathbf{v}\| = 1$ . The plots on the right show results for  $n = 2, 3, 5, 10, 20$ , from above to below. Dashed lines depict the approximation  $\frac{1}{3r^{n-1}}$ .

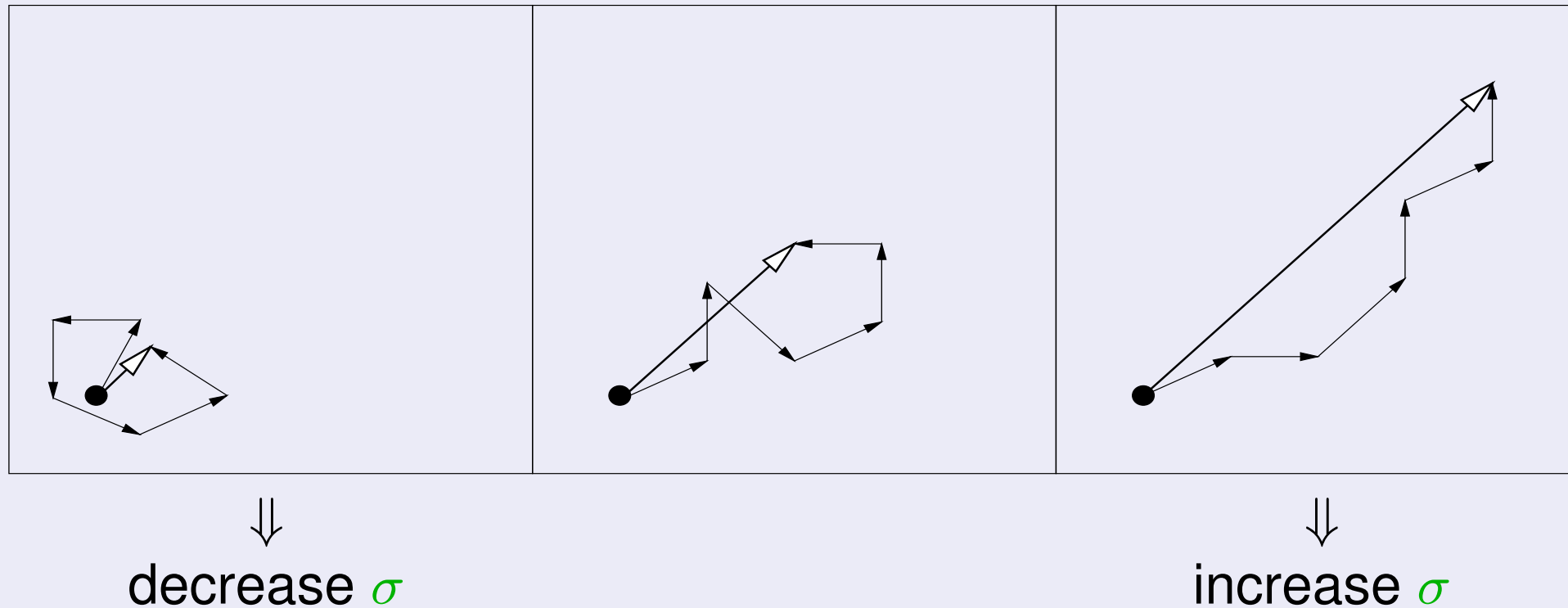
# Path Length Control (CSA)

## 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 generation sequence



loosely speaking steps are

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

# Path Length Control (CSA)

## The Equations

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

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

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

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

# Path Length Control (CSA)

## The Equations

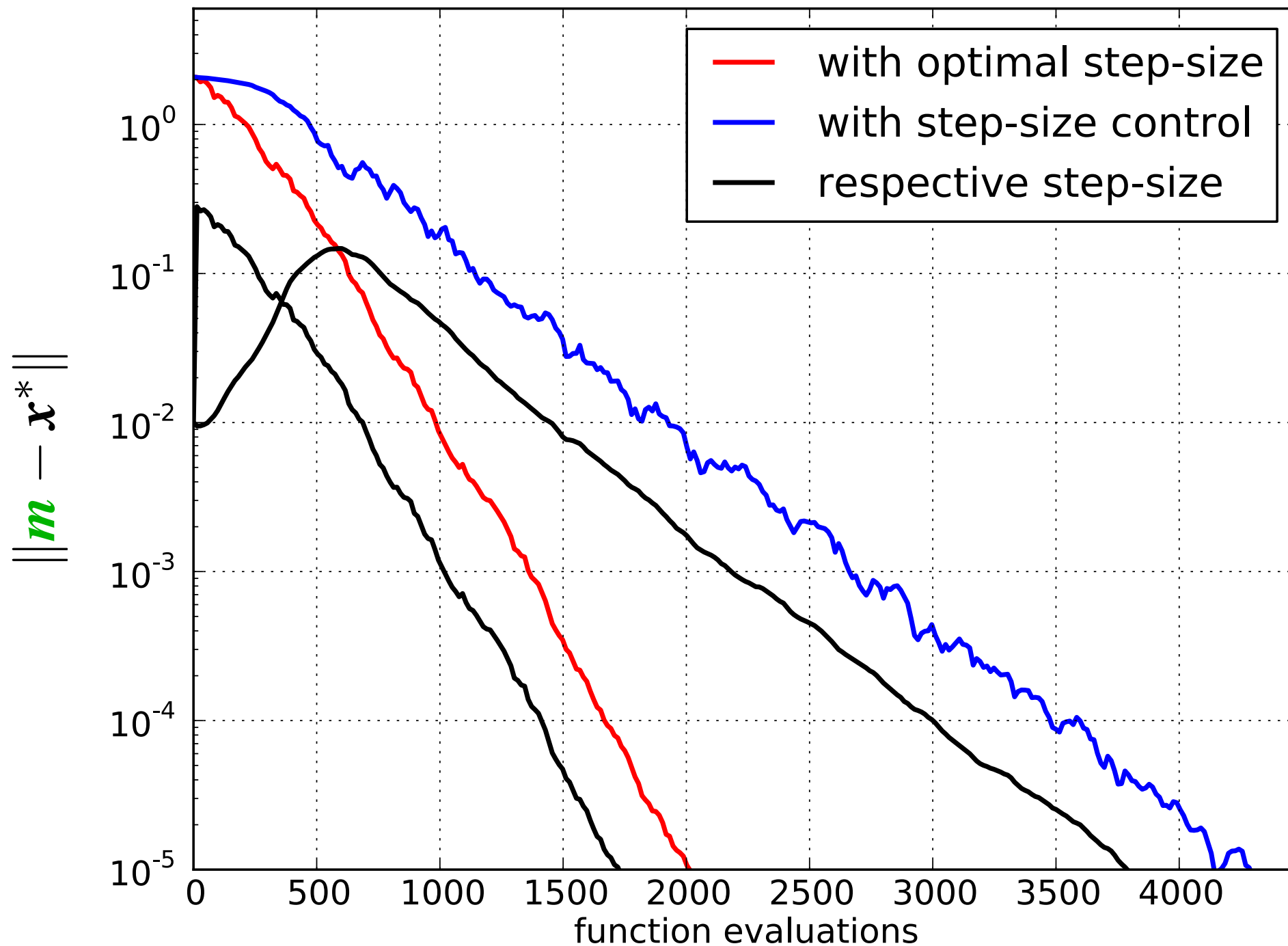
Initialize  $\mathbf{m} \in \mathbb{R}^n$ ,  $\sigma \in \mathbb{R}_+$ , evolution path  $\mathbf{p}_\sigma = \mathbf{0}$ ,  
set  $c_\sigma \approx 4/n$ ,  $d_\sigma \approx 1$ .

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

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

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

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



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

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

# Covariance Matrix Adaptation

# Evolution Strategies

Recalling

New search points are sampled normally distributed

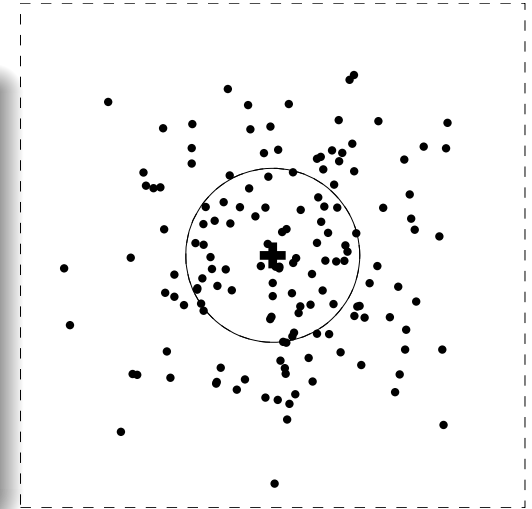
$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

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

where

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

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

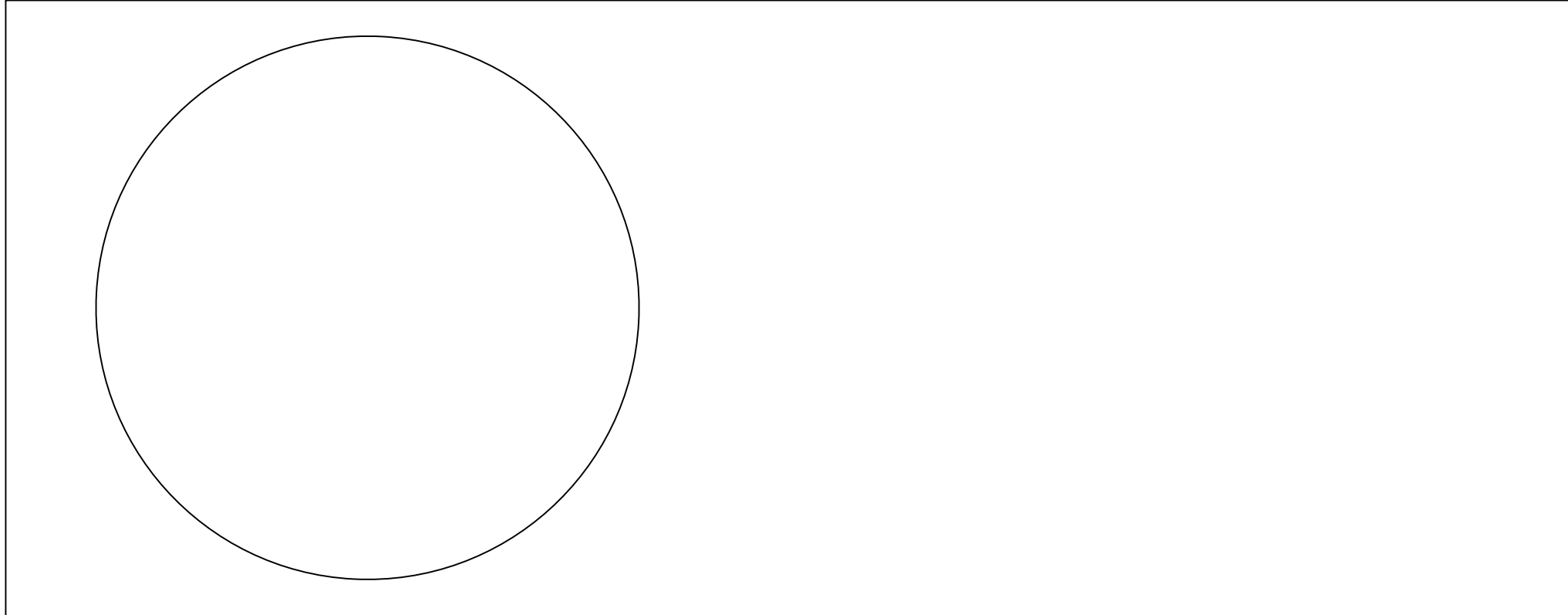




# Covariance Matrix Adaptation

## Rank-One Update

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



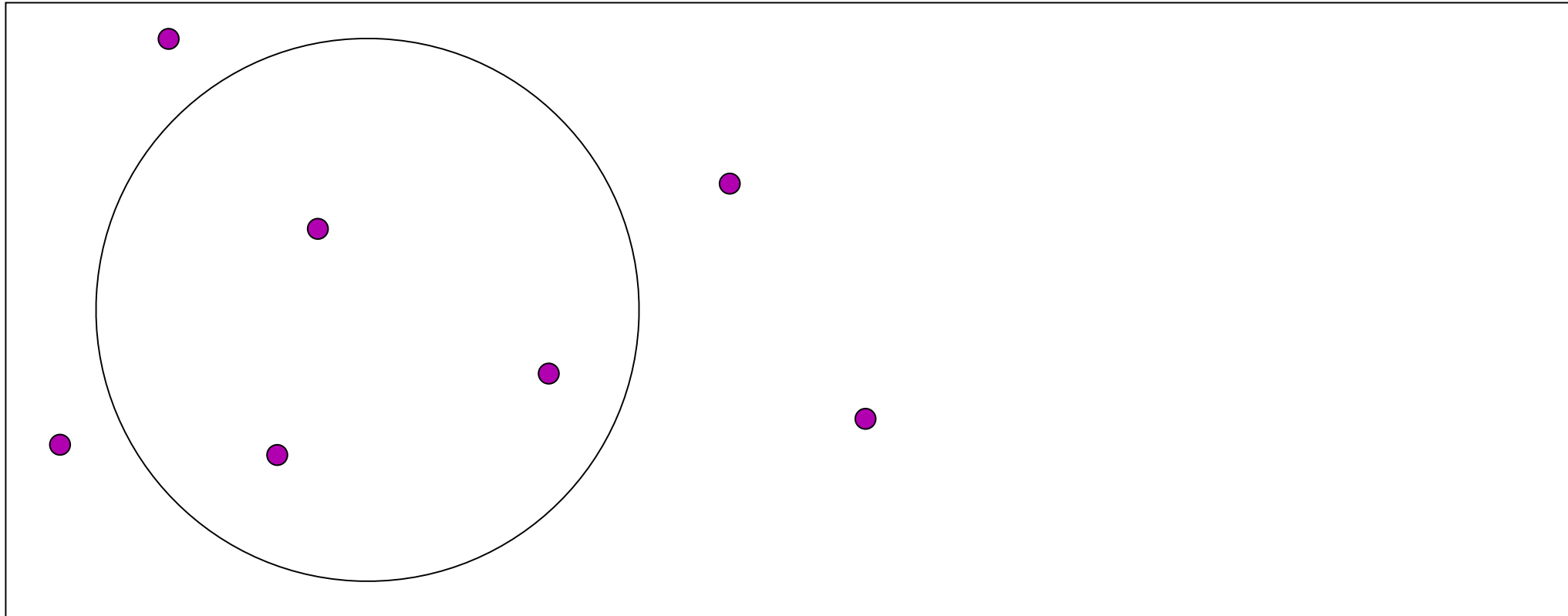
initial distribution,  $\mathbf{C} = \mathbf{I}$

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



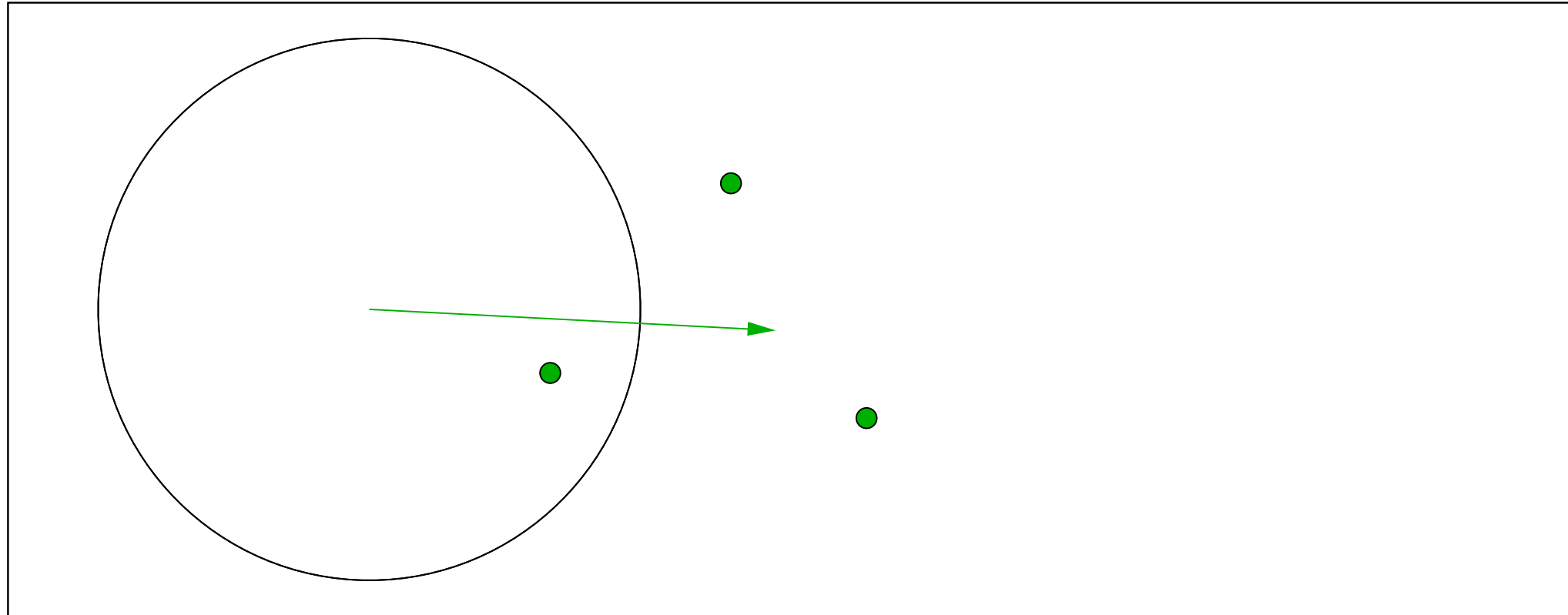
initial distribution,  $\mathbf{C} = \mathbf{I}$

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



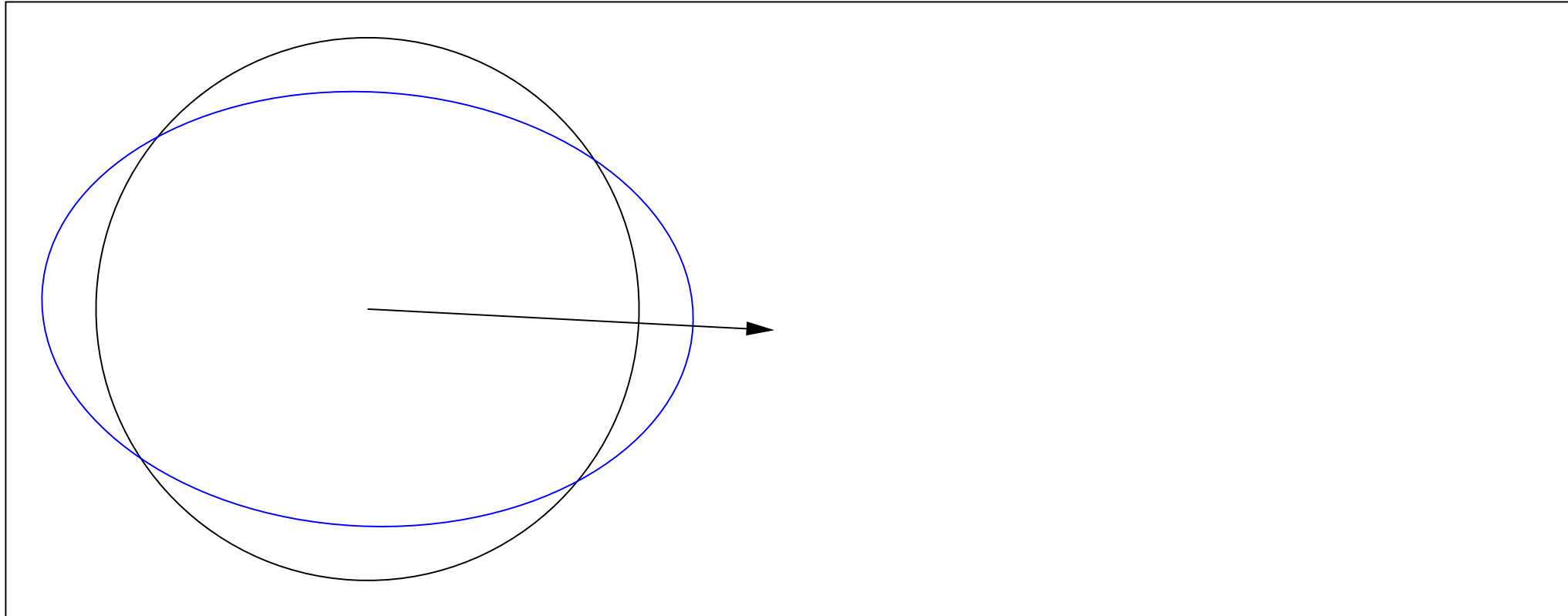
$\mathbf{y}_w$ , movement of the population mean  $\mathbf{m}$  (disregarding  $\sigma$ )

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



mixture of distribution  $\mathbf{C}$  and step  $\mathbf{y}_w$ ,

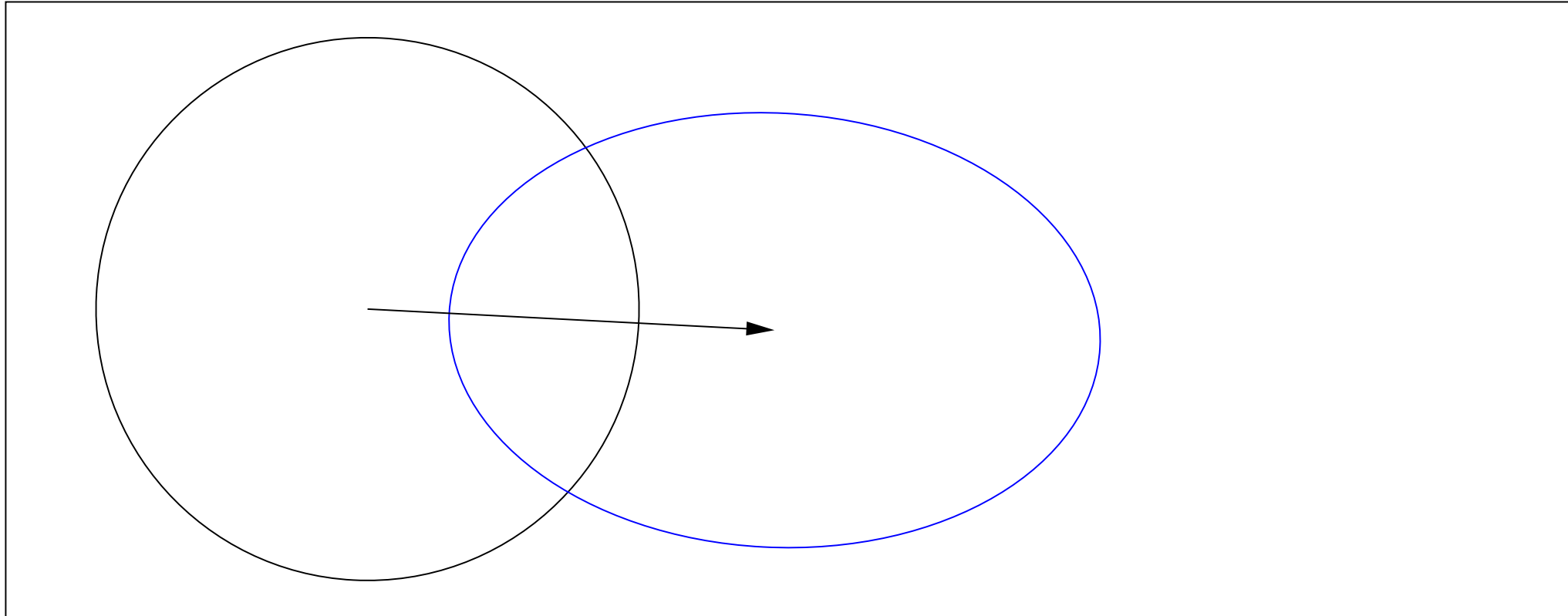
$$\mathbf{C} \leftarrow 0.8 \times \mathbf{C} + 0.2 \times \mathbf{y}_w \mathbf{y}_w^T$$

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



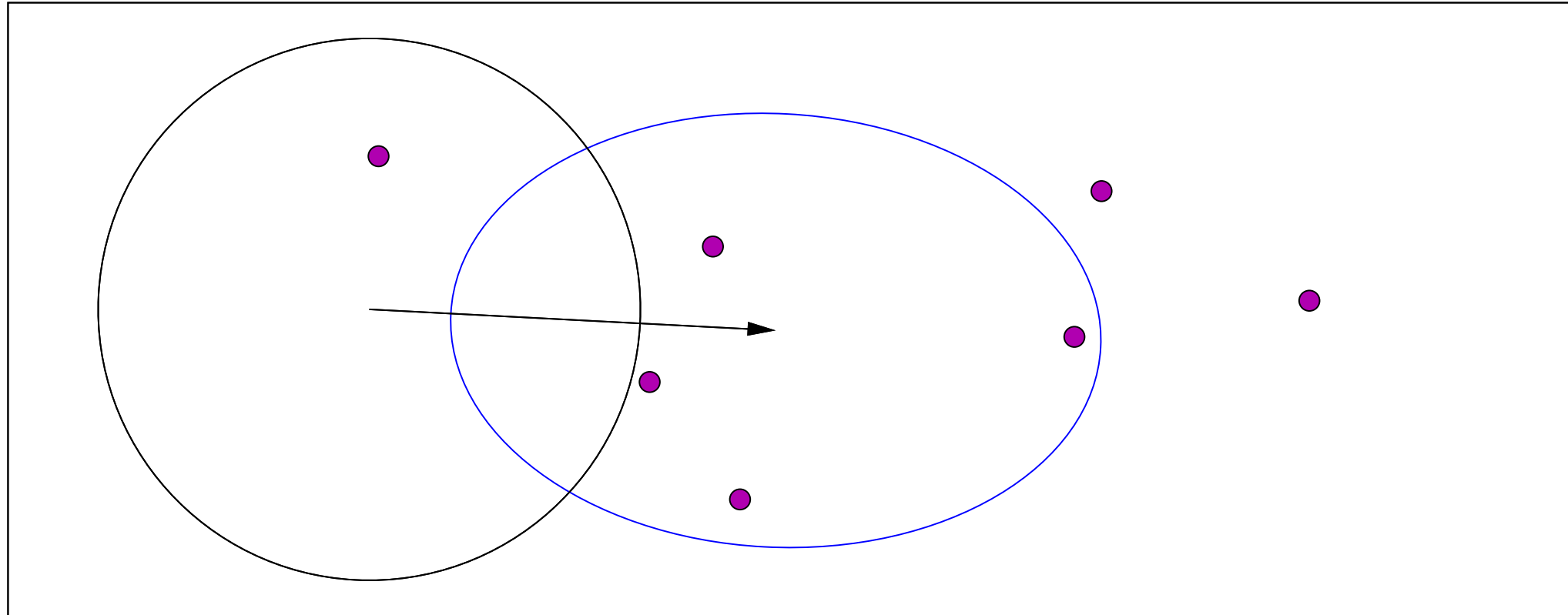
new distribution (disregarding  $\sigma$ )

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



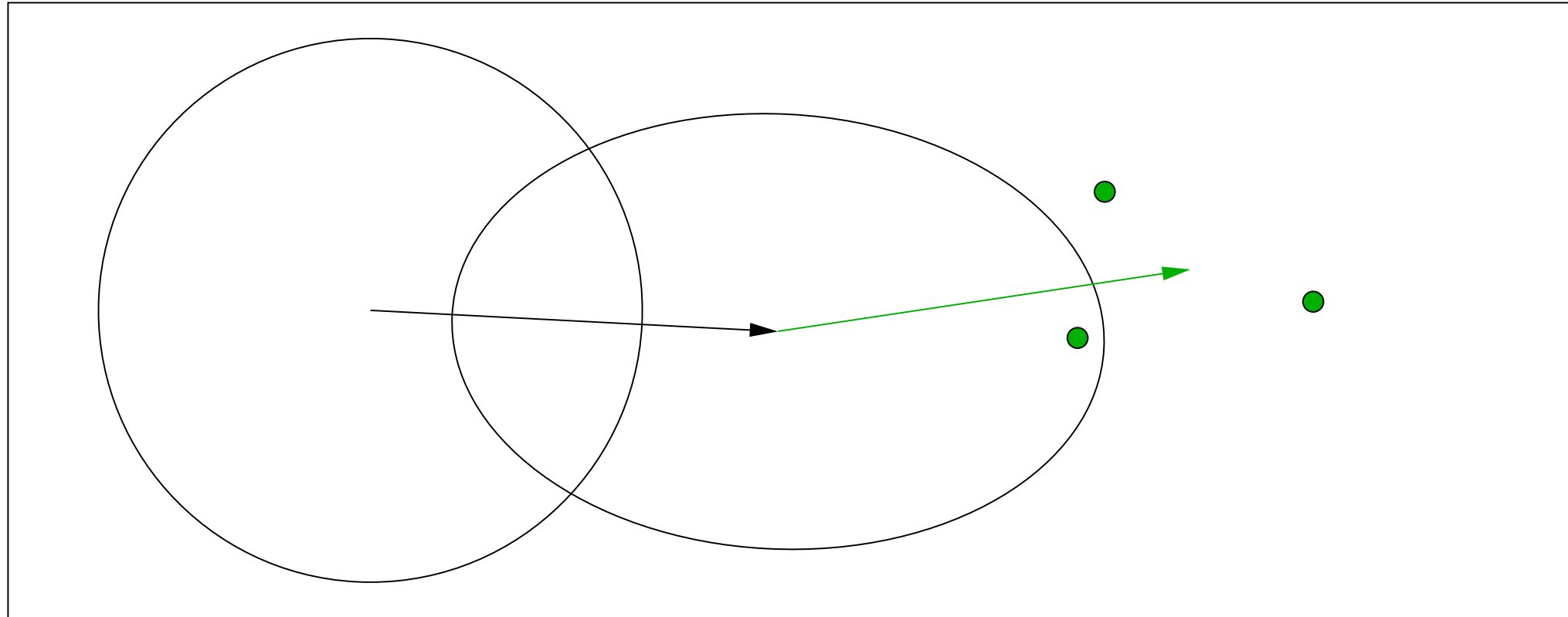
new distribution (disregarding  $\sigma$ )

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



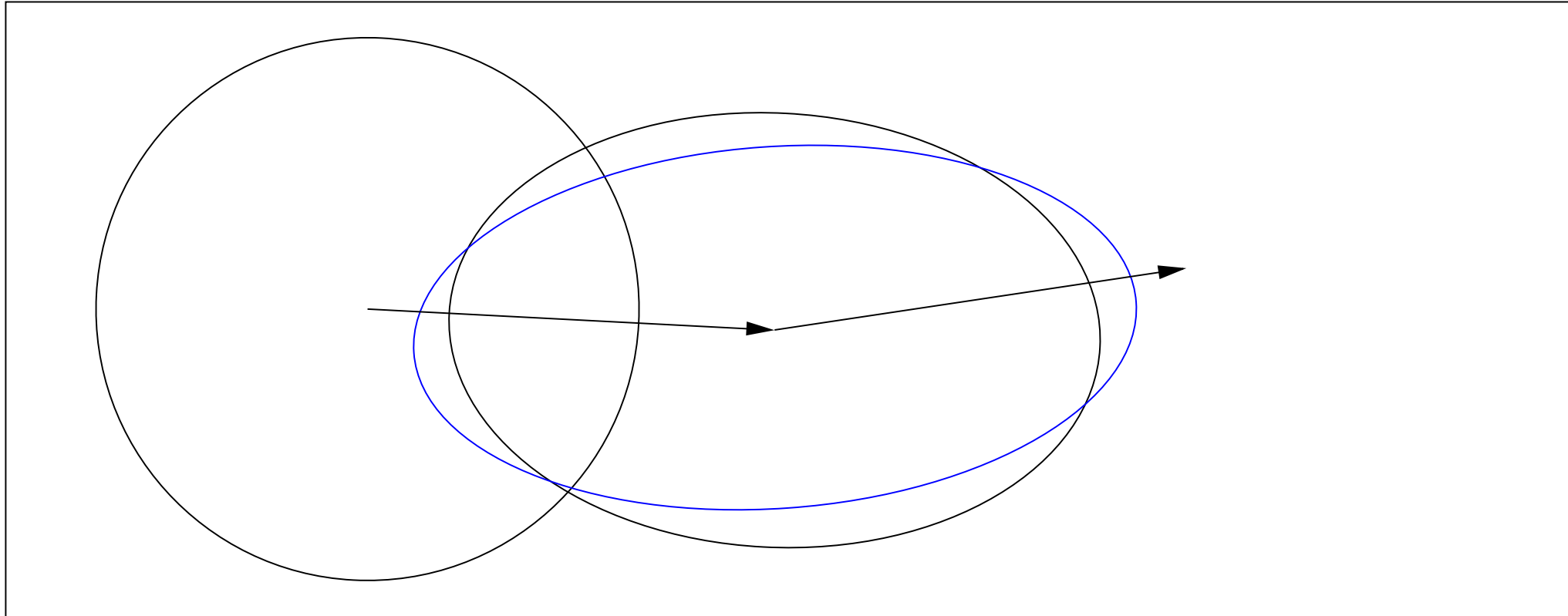
movement of the population mean  $\mathbf{m}$

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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



mixture of distribution  $\mathbf{C}$  and step  $\mathbf{y}_w$ ,

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

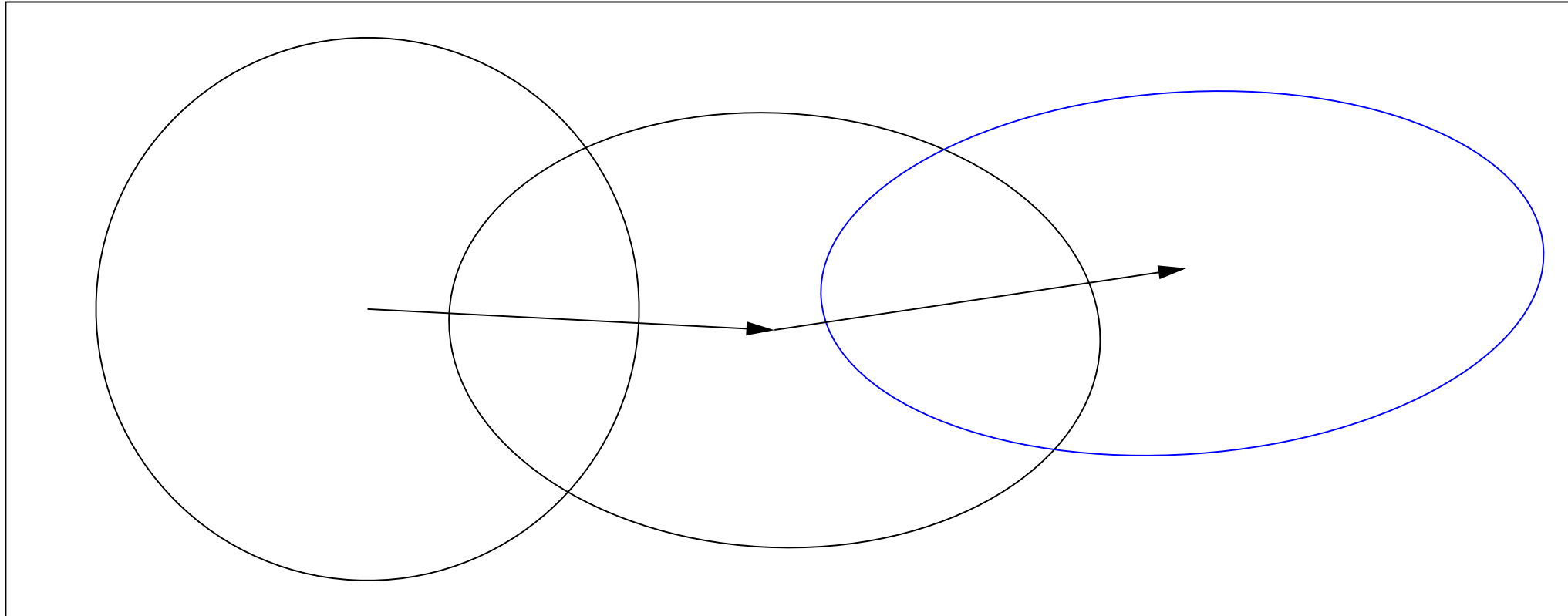
... equations



# Covariance Matrix Adaptation

## Rank-One Update

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



new distribution,

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

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

another viewpoint: the adaptation **follows a natural gradient**

approximation of the expected fitness

... equations

# Covariance Matrix Adaptation

## Rank-One Update

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

While not terminate

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

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

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

The rank-one update has been found independently in several domains<sup>6 7 8 9</sup>

<sup>6</sup> Kjellström&Taxén 1981. Stochastic Optimization in System Design, IEEE TCS

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

<sup>8</sup> Ljung 1999. System Identification: Theory for the User

<sup>9</sup> Haario et al 2001. An adaptive Metropolis algorithm, JSTOR

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

## 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**  
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

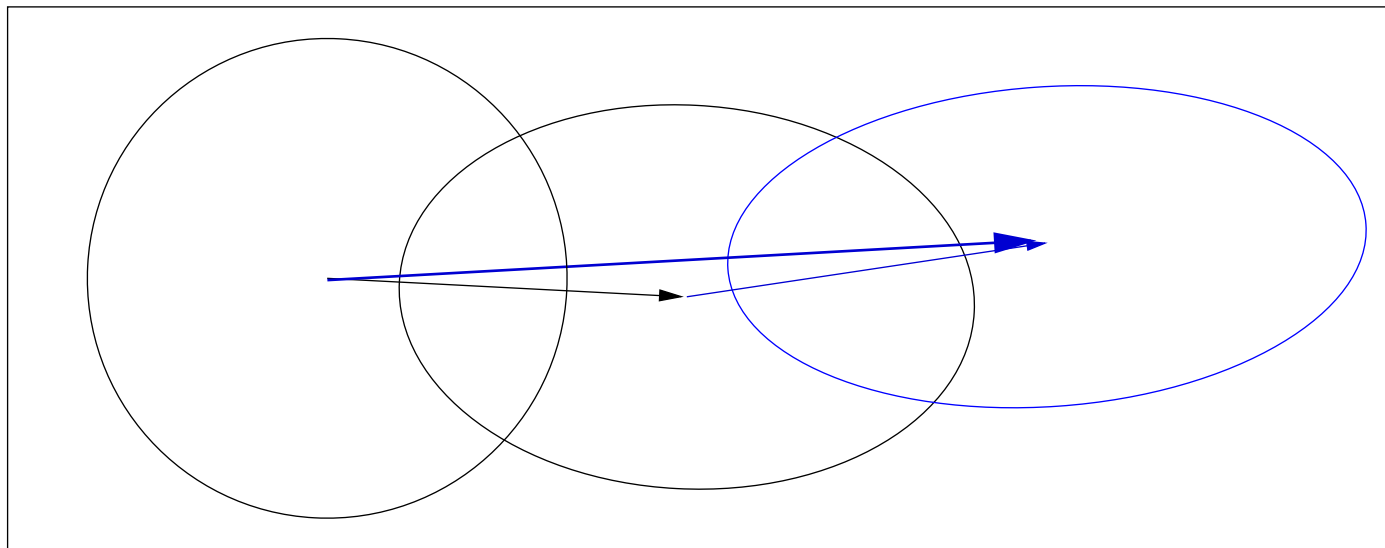
The Evolution Path  
or  
Cumulation

# Cumulation

## The Evolution Path

### Evolution Path

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



An exponentially weighted sum of steps  $y_w$  is used

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

The recursive construction of the evolution path (cumulation):

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

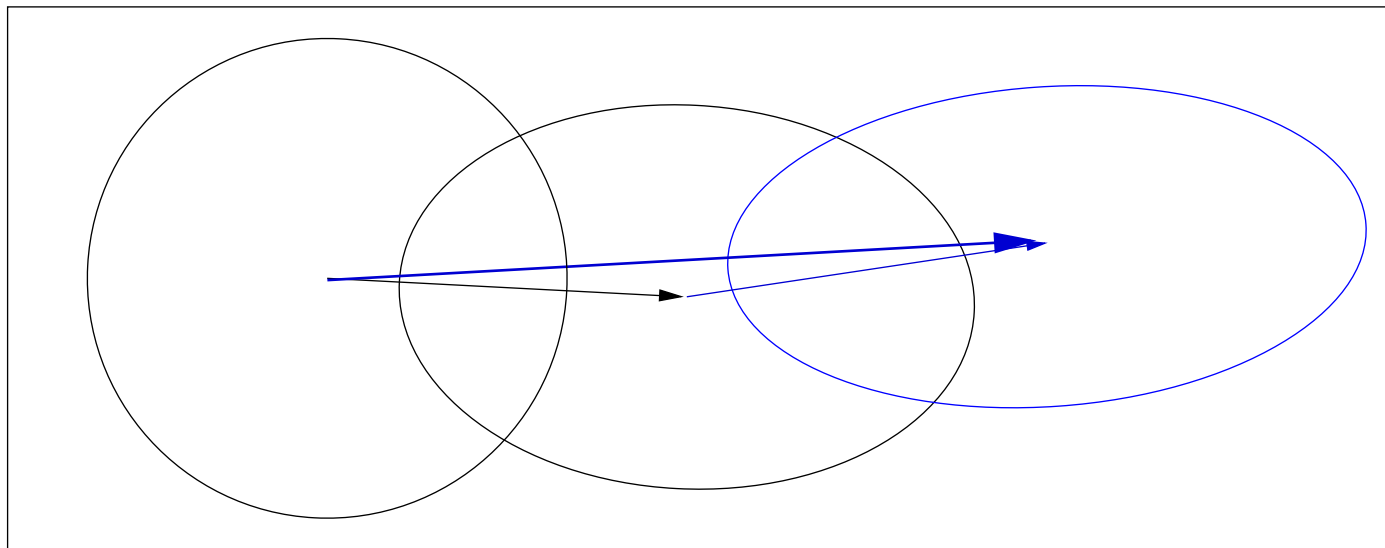
where  $\mu_w = \frac{1}{\sum w_i^2}$ ,  $c_c \ll 1$ . **History information** is accumulated in the evolution path.

# Cumulation

## The Evolution Path

### Evolution Path

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



An exponentially weighted sum of steps  $y_w$  is used

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

The recursive construction of the evolution path (cumulation):

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

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

“Cumulation” is a widely used technique and also know as

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

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

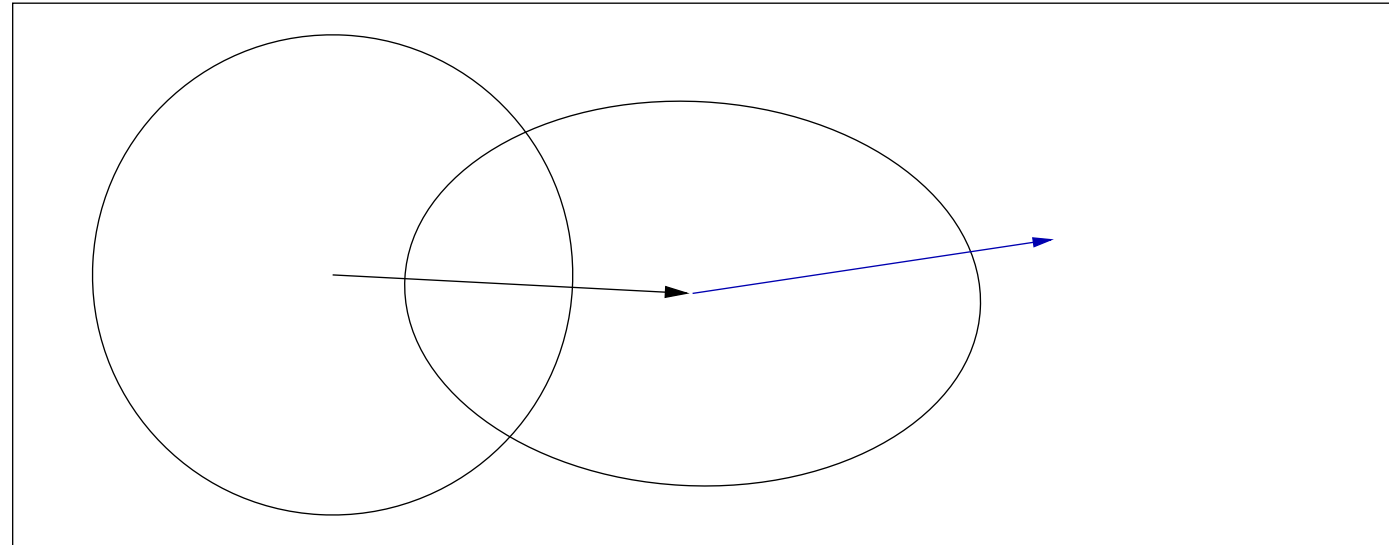
...why?

# Cumulation

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

## Utilizing the Evolution Path

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



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

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

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

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

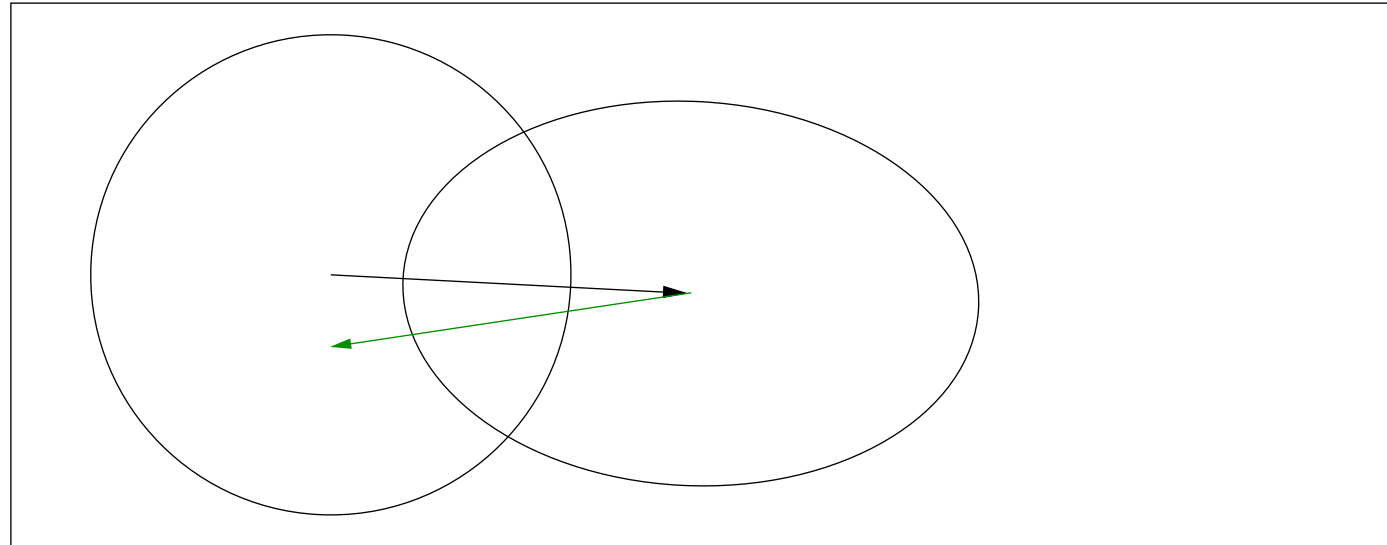


# Cumulation

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

## Utilizing the Evolution Path

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



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

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

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

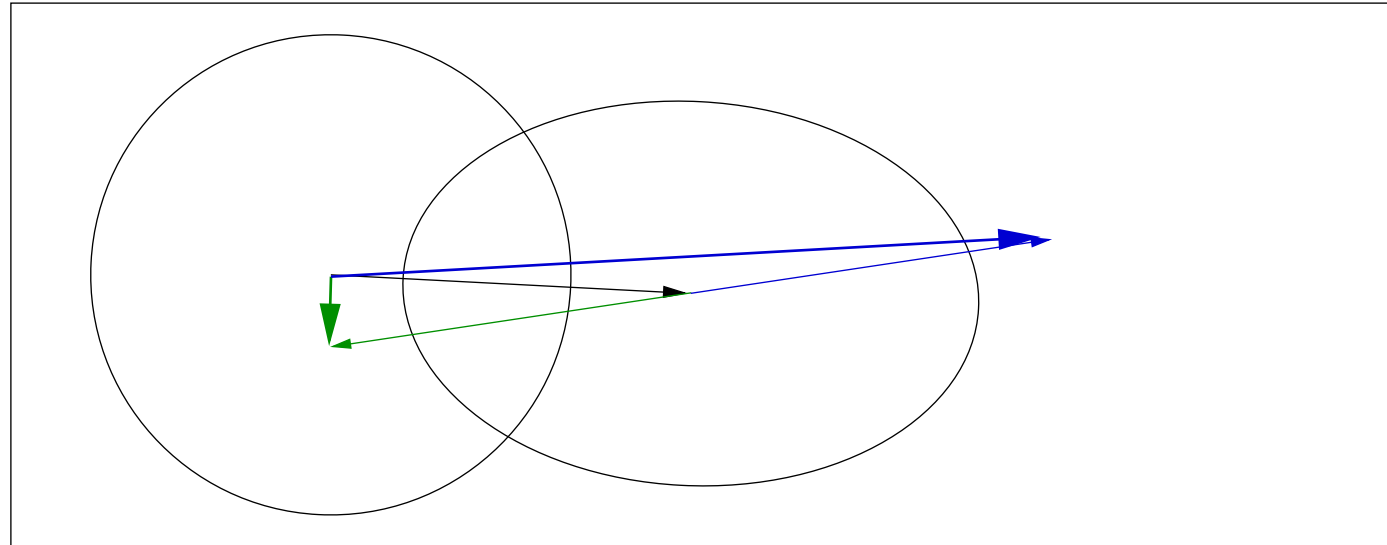
where  $\mu_w = \frac{1}{\sum w_i^2}$ ,  $c_{\text{cov}} \ll c_c \ll 1$  such that  $1/c_c$  is the “backward time horizon”.

# Cumulation

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

## Utilizing the Evolution Path

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



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

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

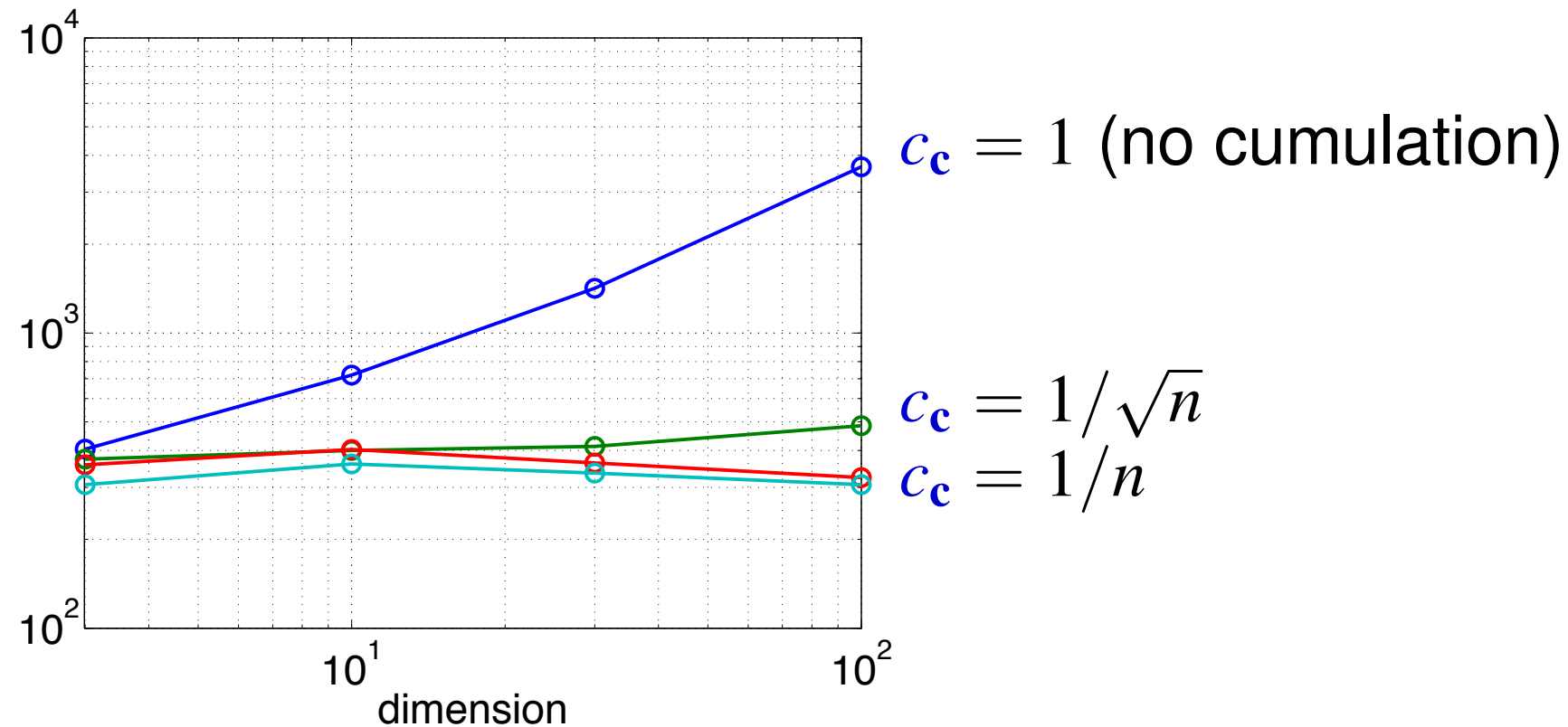
$$\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + c_{\text{cov}} \underbrace{\mathbf{p}_c\mathbf{p}_c^T}_{\text{rank-one}}$$

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

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  $\mathcal{O}(n)$** .<sup>(a)</sup>

<sup>a</sup>Hansen & Auger 2013. Principled design of continuous stochastic search: From theory to practice.

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



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

# Rank- $\mu$ 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- $\mu$  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

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

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 <sup>10</sup>

The rank- $\mu$  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  $c_{\text{cov}} \leq 1$ .

<sup>10</sup>Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. 

# Rank- $\mu$ 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- $\mu$  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

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

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 <sup>10</sup>

The rank- $\mu$  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  $c_{\text{cov}} \leq 1$ .

<sup>10</sup>Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.

# Rank- $\mu$ 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- $\mu$  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

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

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

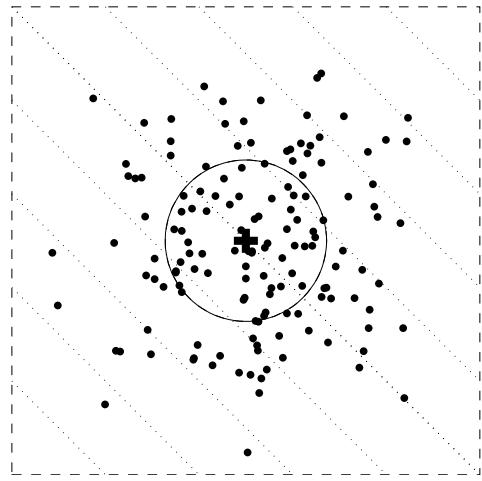
with  $\mu = \lambda$  weights can be negative <sup>10</sup>

The rank- $\mu$  update then reads

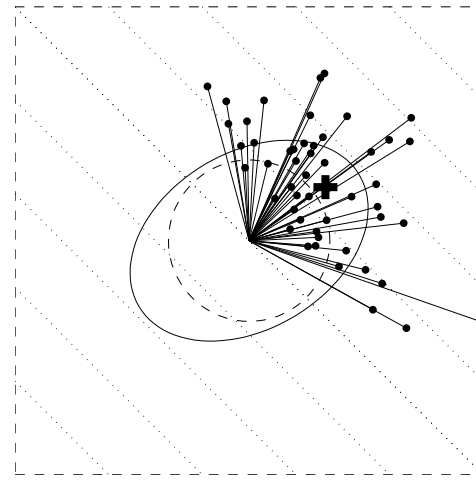
$$\mathbf{C} \leftarrow (1 - c_{\text{cov}}) \mathbf{C} + c_{\text{cov}} \mathbf{C}_\mu$$

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

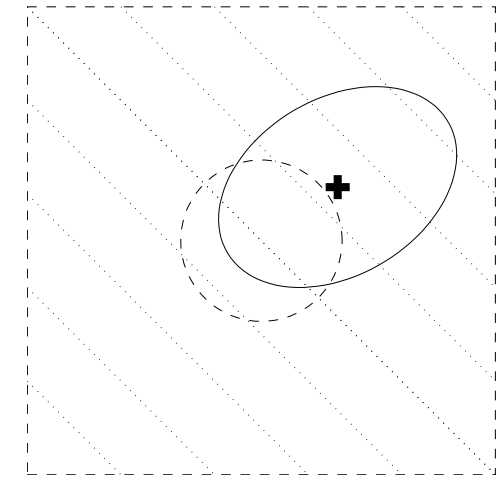
<sup>10</sup>Jastrebski and Arnold (2006). Improving evolution strategies through active covariance matrix adaptation. CEC.



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



$$\begin{aligned} \mathbf{C}_\mu &= \frac{1}{\mu} \sum y_{i:\lambda} y_{i:\lambda}^\top \\ \mathbf{C} &\leftarrow (1 - 1) \times \mathbf{C} + 1 \times \mathbf{C}_\mu \end{aligned}$$

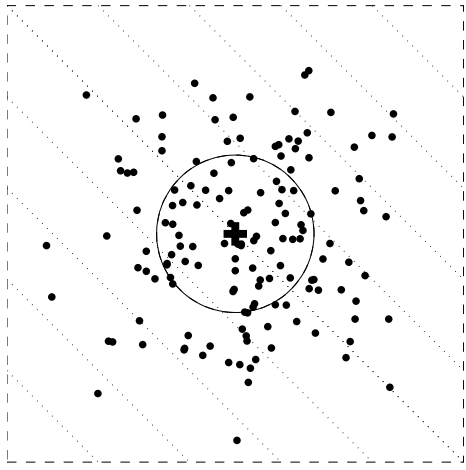


$$m_{\text{new}} \leftarrow m + \frac{1}{\mu} \sum y_{i:\lambda}$$

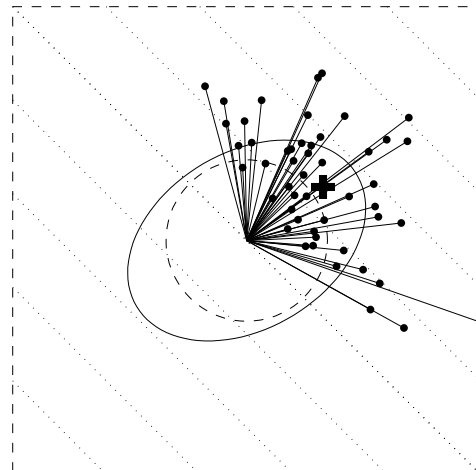
new distribution

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

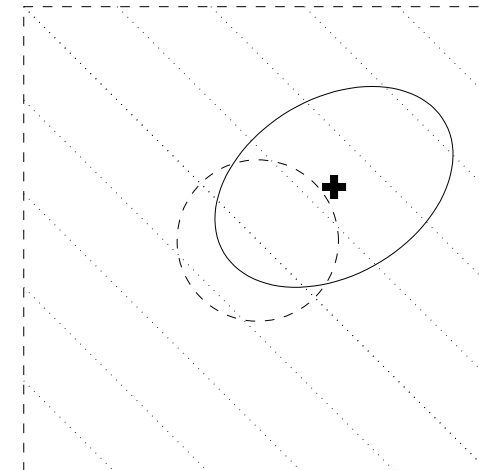
calculating  $\mathbf{C}$  where  
 $\mu = 50$ ,  
 $w_1 = \dots = w_\mu = \frac{1}{\mu}$ ,  
and  $c_{\text{cov}} = 1$

Rank- $\mu$  CMA versus Estimation of Multivariate Normal Algorithm EMNA<sub>global</sub><sup>11</sup>

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

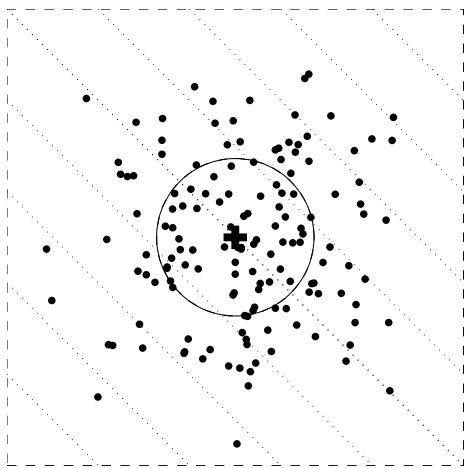


$$\mathbf{C} \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - m_{\text{old}})(x_{i:\lambda} - m_{\text{old}})^T$$

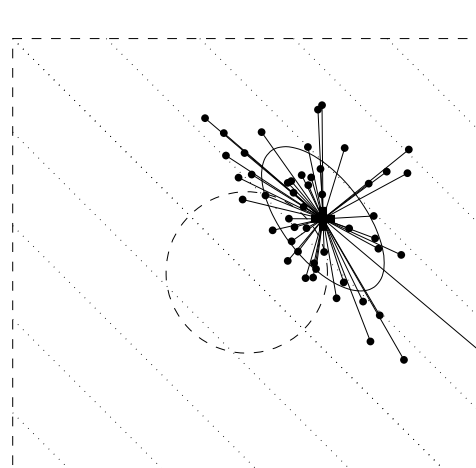


$$m_{\text{new}} = m_{\text{old}} + \frac{1}{\mu} \sum y_{i:\lambda}$$

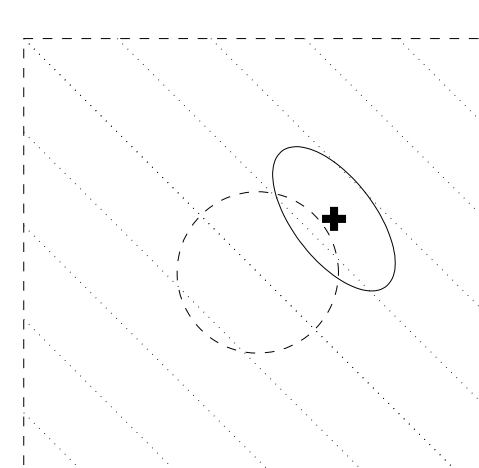
rank- $\mu$  CMA  
conducts a  
PCA of  
steps



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



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



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

EMNA<sub>global</sub>  
conducts a  
PCA of  
points

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

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

new distribution

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

<sup>11</sup> 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  $\mathcal{O}(n)$  <sup>(12)</sup>  
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

<sup>12</sup>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  $\mathcal{O}(n)$  <sup>(12)</sup>  
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

<sup>12</sup>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  $\mathcal{O}(n)$  <sup>(12)</sup>  
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

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

# Summary of Equations

## The Covariance Matrix Adaptation Evolution Strategy

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

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

**Set:**  $c_c \approx 4/n$ ,  $c_\sigma \approx 4/n$ ,  $c_1 \approx 2/n^2$ ,  $c_\mu \approx \mu_w/n^2$ ,  $c_1 + c_\mu \leq 1$ ,  $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$ ,  
and  $w_{i=1\dots\lambda}$  such that  $\mu_w = \frac{1}{\sum_{i=1}^{\mu} w_i^2} \approx 0.3 \lambda$

**While not terminate**

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

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

$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbb{1}_{\{\|\mathbf{p}_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w$  cumulation for  $\mathbf{C}$

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$  update  $\mathbf{C}$

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

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

# Source Code Snippet

```

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

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

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

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

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

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

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

```

# Strategy Internal Parameters

- related to selection and recombination
  - ▶  $\lambda$ , offspring number, new solutions sampled, population size
  - ▶  $\mu$ , parent number, solutions involved in updates of  $m$ ,  $C$ , and  $\sigma$
  - ▶  $w_{i=1,\dots,\mu}$ , recombination weights
- related to  $C$ -update
  - ▶  $c_c$ , decay rate for the evolution path
  - ▶  $c_1$ , 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

Parameters were identified in carefully chosen experimental set ups. **Parameters do not in the first place depend on the objective function** and are not meant to be in the users choice.

Only(?) the population size  $\lambda$  (and the initial  $\sigma$ ) might be reasonably varied in a wide range, *depending on the objective function*

Useful: restarts with increasing population size (IPOP)

# Experimentum Crucis (0)

What did we want to achieve?

- reduce any convex-quadratic function

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

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

to the sphere model

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

without use of derivatives

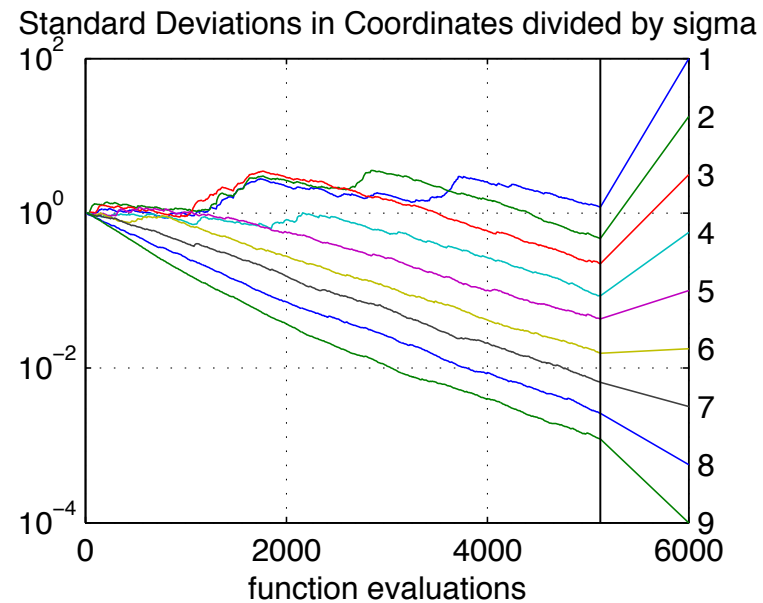
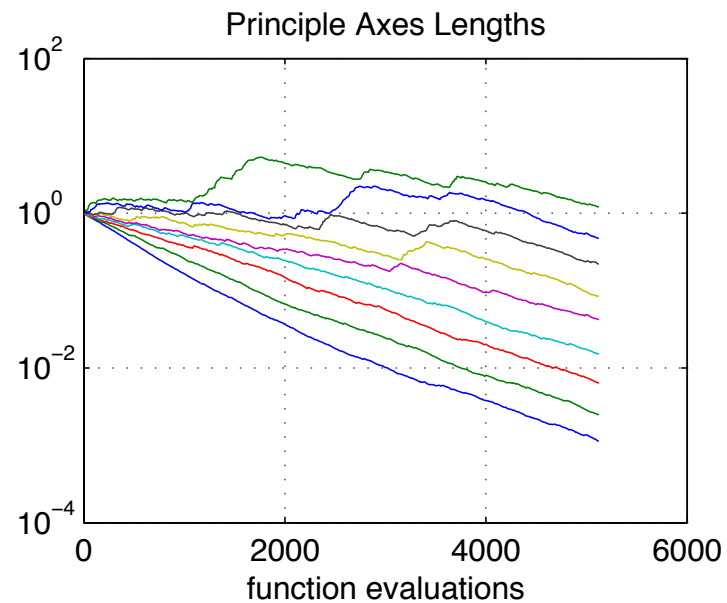
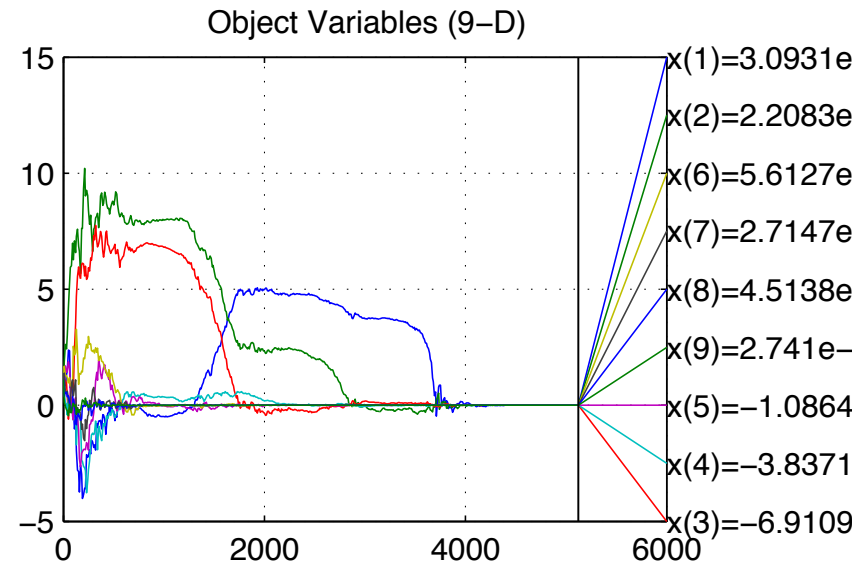
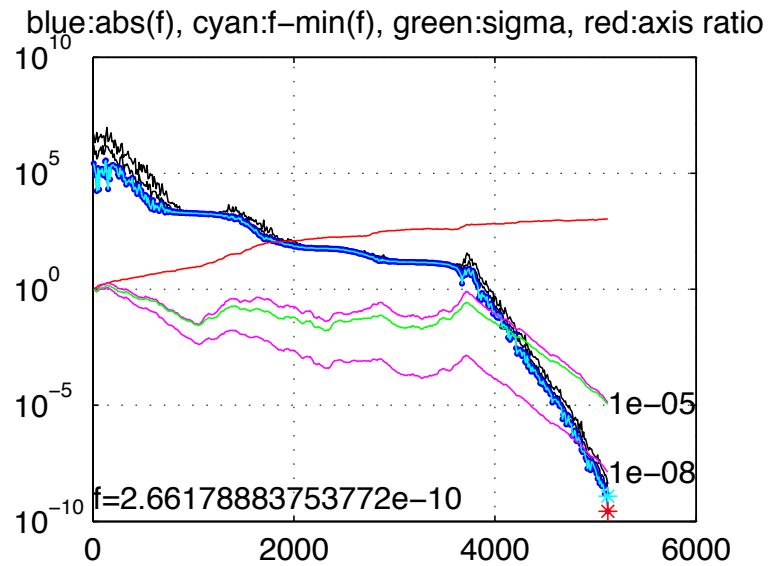
- lines of equal density align with lines of equal fitness

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

in a stochastic sense

# Experimentum Crucis (1)

$f$  convex quadratic, separable

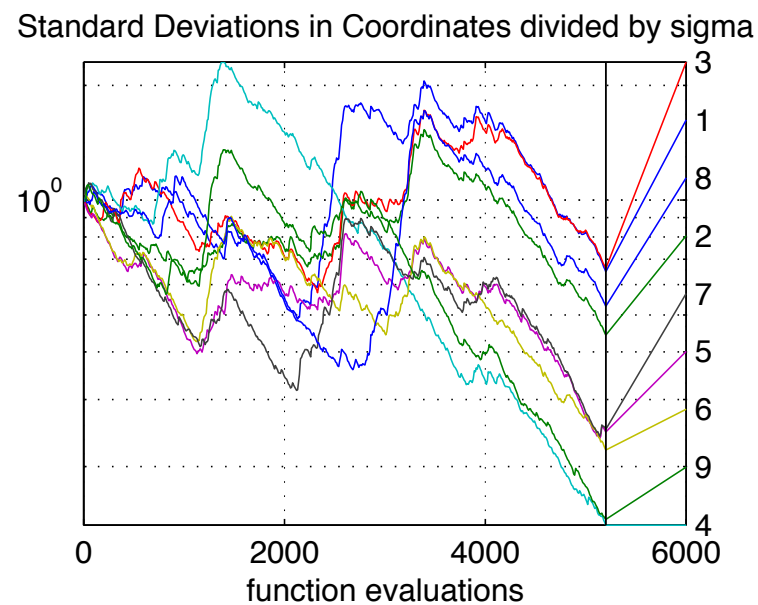
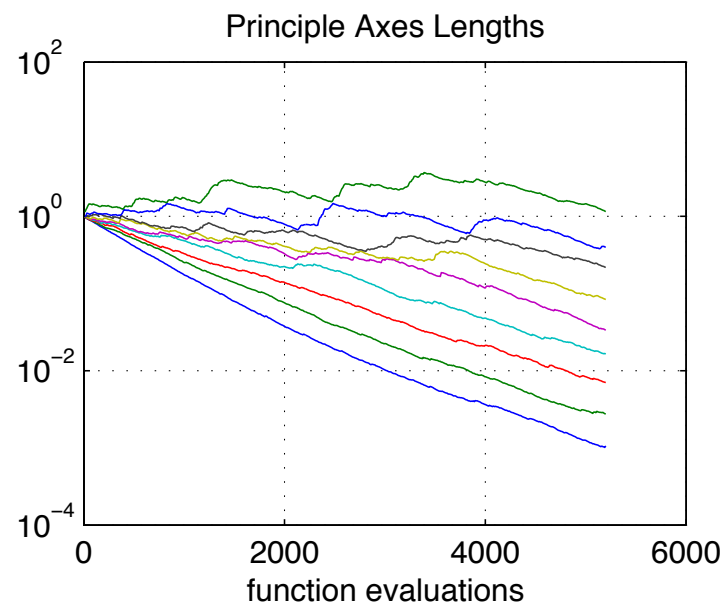
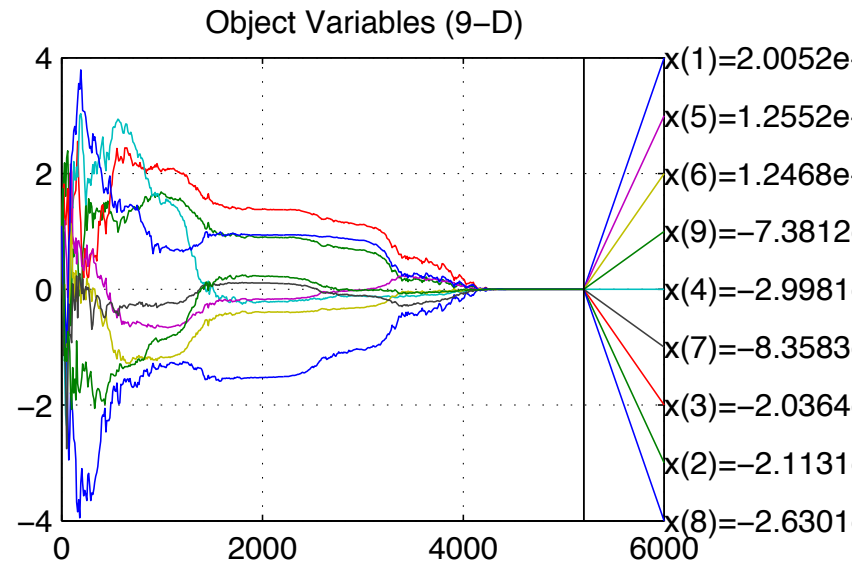
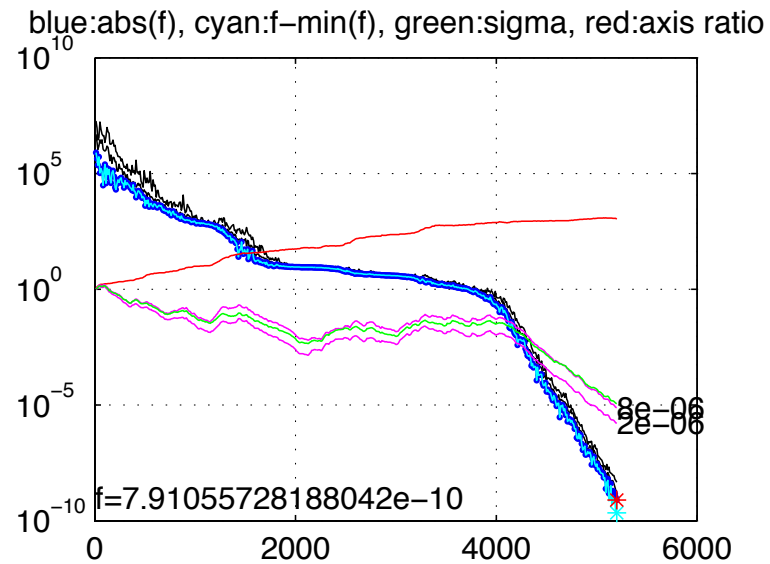


$$f(\mathbf{x}) = \sum_{i=1}^n 10^{\alpha \frac{i-1}{n-1}} x_i^2, \alpha = 6$$



# Experimentum Crucis (2)

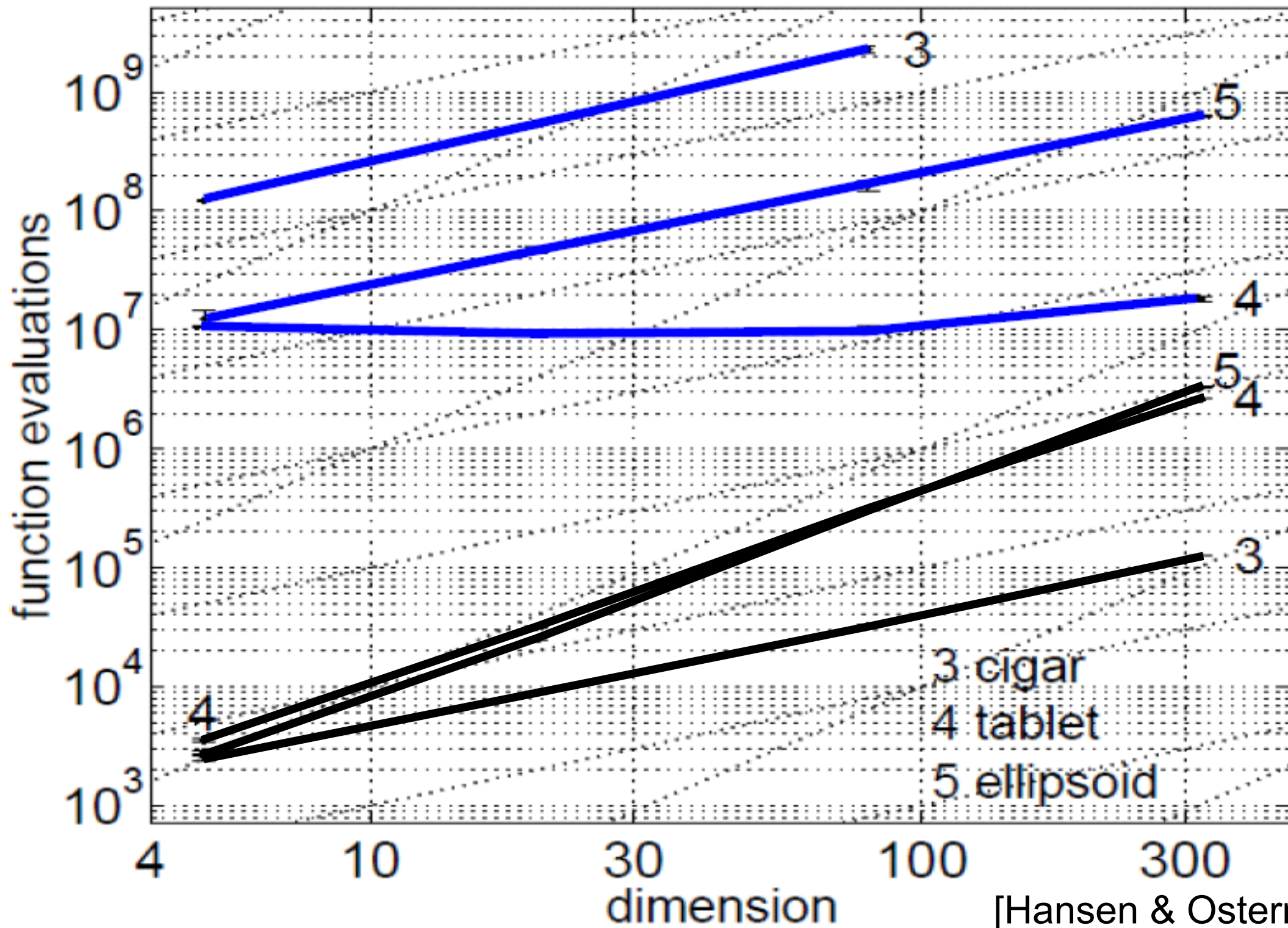
$f$  convex quadratic, as before but non-separable (rotated)



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

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

# Quantifying the enhancement



[Hansen & Ostermeier 2001]

black: CMA-ES ( $c_1 \approx 2/n^2$ ), blue: CSA-ES ( $c_1 = 0$ )

# Theoretical Considerations

# CMA-ES = Natural Evolution Strategy + Cumulation

Natural gradient descend using the MC approximation and the normal distribution

- Rewriting the update of the distribution mean

$$\mathbf{m}_{\text{new}} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \underbrace{\sum_{i=1}^{\mu} w_i (\mathbf{x}_{i:\lambda} - \mathbf{m})}_{\text{natural gradient for mean } \frac{\tilde{\partial}}{\partial \mathbf{m}} \hat{\mathbb{E}}(w \circ P_f(f(\mathbf{x})) | \mathbf{m}, \mathbf{C})}$$

- Rewriting the update of the covariance matrix<sup>13</sup>

$$\begin{aligned} \mathbf{C}_{\text{new}} \leftarrow & \mathbf{C} + c_1 \overbrace{(\mathbf{p}_c \mathbf{p}_c^T - \mathbf{C})}^{\text{rank one}} \\ & + \underbrace{\frac{c_\mu}{\sigma^2} \sum_{i=1}^{\mu} w_i \left( \overbrace{(\mathbf{x}_{i:\lambda} - \mathbf{m})(\mathbf{x}_{i:\lambda} - \mathbf{m})^T}_{\text{rank-}\mu} - \sigma^2 \mathbf{C} \right)}_{\text{natural gradient for covariance matrix } \frac{\tilde{\partial}}{\partial \mathbf{C}} \hat{\mathbb{E}}(w \circ P_f(f(\mathbf{x})) | \mathbf{m}, \mathbf{C})} \end{aligned}$$

<sup>13</sup> Akimoto et.al. (2010): Bidirectional Relation between CMA Evolution Strategies and Natural Evolution Strategies, PPSN XI

# Summary of Equations

The Covariance Matrix Adaptation Evolution Strategy

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

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

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

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

**While not terminate**

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

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

$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbb{1}_{\{\|\mathbf{p}_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w$  cumulation for  $\mathbf{C}$

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

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$  update  $\mathbf{C}$

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

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

# Maximum Likelihood Update

The new distribution mean  $\mathbf{m}$  maximizes the log-likelihood

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

independently of the given covariance matrix

$$\log p_{\mathcal{N}}(\mathbf{x} | \mathbf{m}, \mathbf{C}) = -\frac{1}{2} \log \det(2\pi \mathbf{C}) - \frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})$$

$p_{\mathcal{N}}$  is the density of the multi-variate normal distribution

# Maximum Likelihood Update

The new distribution mean  $\mathbf{m}$  maximizes the log-likelihood

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

independently of the given covariance matrix

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

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

$$\log p_{\mathcal{N}}(\mathbf{x} | \mathbf{m}, \mathbf{C}) = -\frac{1}{2} \log \det(2\pi \mathbf{C}) - \frac{1}{2} (\mathbf{x} - \mathbf{m})^{\text{T}} \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})$$

$p_{\mathcal{N}}$  is the density of the multi-variate normal distribution



# Variable Metric

On the function class

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

the covariance matrix approximates the inverse Hessian up to a constant factor, that is:

$$\mathbf{C} \propto \mathbf{H}^{-1} \quad (\text{approximately})$$

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

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



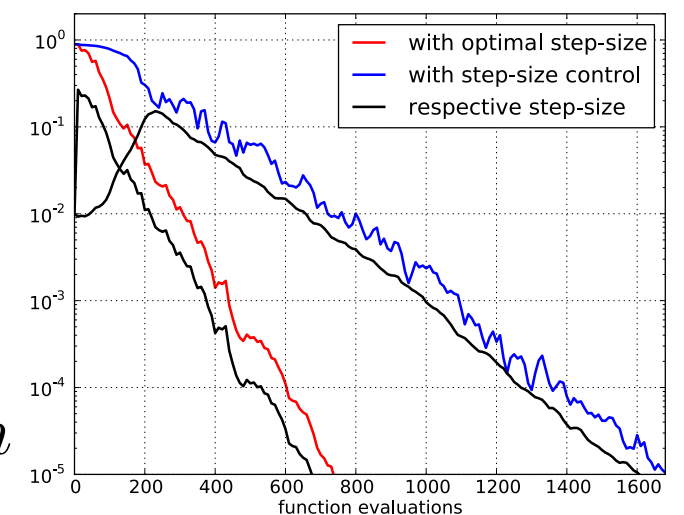
# On Convergence

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

$$\|\mathbf{m}_k - \mathbf{x}^*\| \propto e^{-ck}, \quad c \leq \frac{0.25}{n}$$

where  $k$  is the number of  $f$ -evaluations. In practice:  $c \approx 0.1/n$

selection	$n \cdot c_{\max}$
(1+1)	0.202
$(\mu/\mu, \lambda)$	0.202
$(\mu/\mu_w, \lambda)$	0.25



Monte Carlo pure random search converges like

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

# On the Sphere Function (for $n$ large)

optimal population size  $\lambda$  (offspring)

$$\lambda \not\ll 5 \text{ and } \lambda \not\gg n \quad \text{for } \mu = 1 \text{ we have } \lambda^{\text{opt}} = 5$$

optimal recombination weights

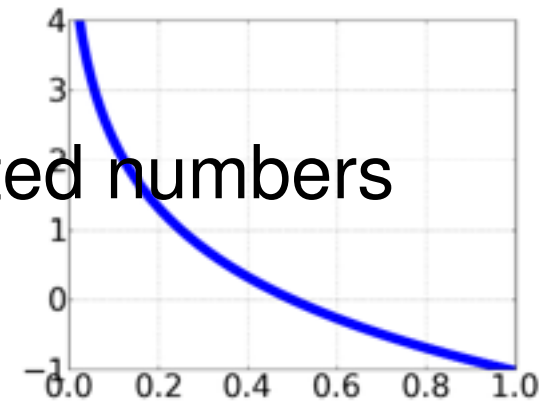
$$w_i^{\text{opt}} \propto -\mathcal{N}_{i:\lambda}(0, 1) \quad i\text{th order statistic of } \lambda \text{ normally distributed numbers}$$

optimal (effective) parent number

$$\mu_w := \frac{1}{\sum_i w_i^2} \leq \mu \quad \text{with } \sum_i |w_i| = 1$$

$$\mu_w^{\text{opt}} \approx 0.32\lambda \quad \text{if } w_i = \max(w_i^{\text{opt}}, 0)$$

$$\mu^{\text{opt}} \approx 0.27\lambda \quad \text{if } w_{i=1\dots\mu} = 1/\mu \text{ (truncation selection)}$$



# On the Sphere Function (for $n$ large)

optimal step-size if  $\mu_w < n$

$$\sigma^{\text{opt}} \propto c_w \frac{\mu_w}{n} \quad \text{where } c_w = - \sum_i w_i \mathcal{N}_{i:\lambda}(0, 1) \approx 1$$

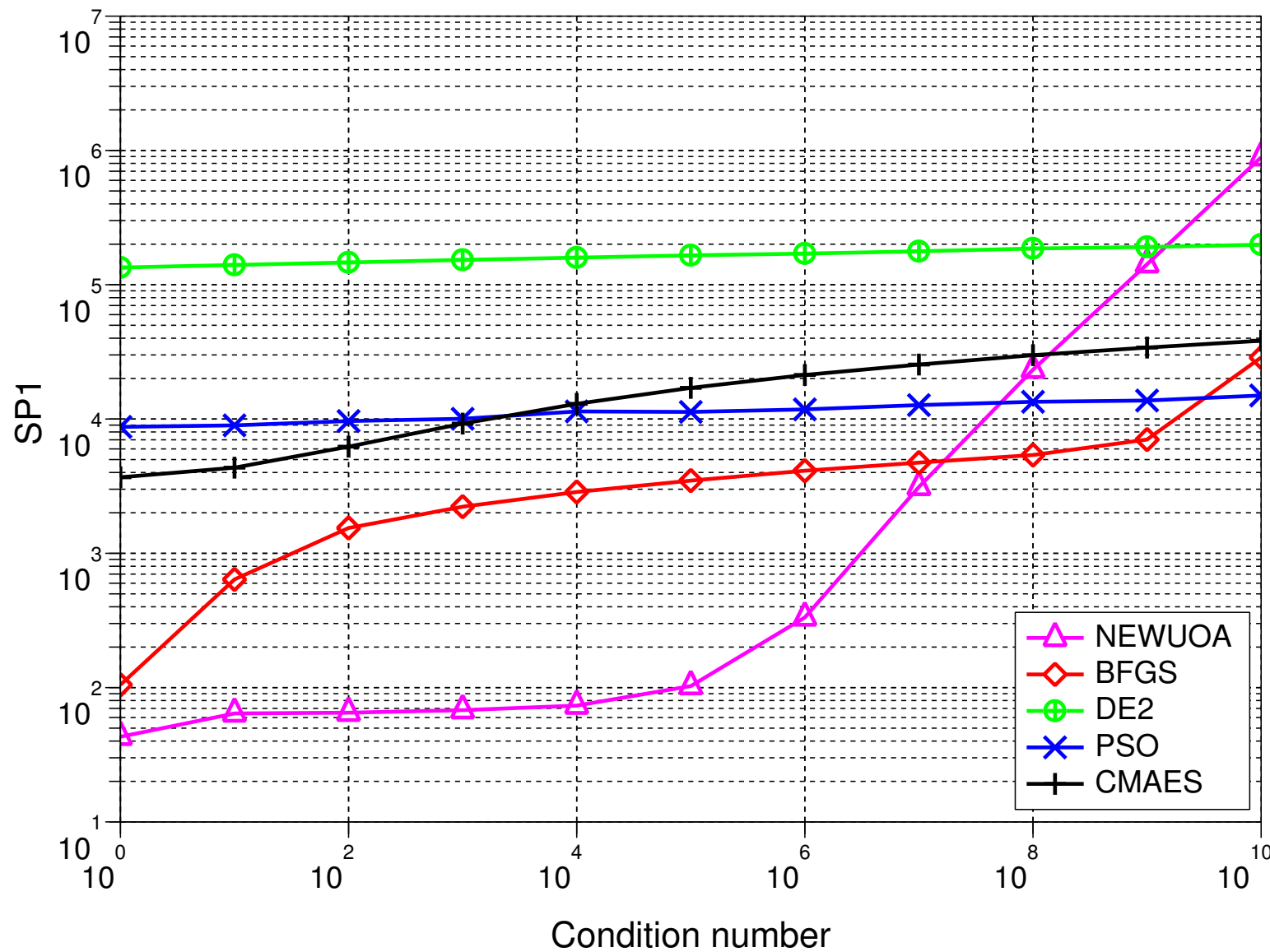
is proportional to  $\mu_w$  and therefore to the population size

# Comparing Experiments

# Comparison to BFGS, NEWUOA, PSO and DE

$f$  convex quadratic, separable with varying condition number  $\alpha$

Ellipsoid dimension 20, 21 trials, tolerance  $1e-09$ , eval max  $1e+07$



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

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

$\mathbf{H}$  diagonal

$g$  identity (for **BFGS** and **NEWUOA**)

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

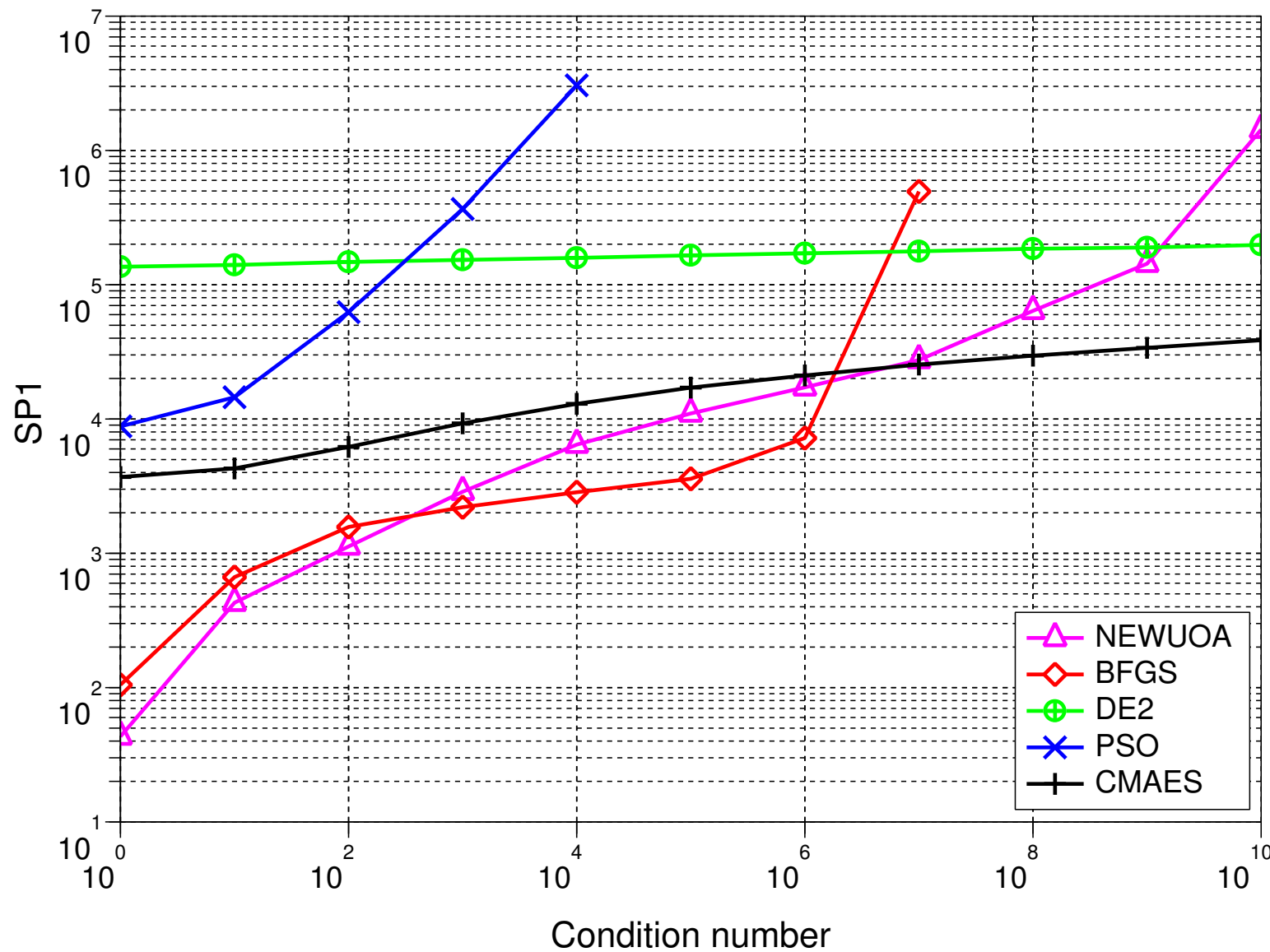
SP1 = average number of objective function evaluations<sup>14</sup> to reach the target function value of  $g^{-1}(10^{-9})$

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

# Comparison to BFGS, NEWUOA, PSO and DE

$f$  convex quadratic, non-separable (rotated) with varying condition number  $\alpha$

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



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

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

$\mathbf{H}$  full

$g$  identity (for **BFGS** and **NEWUOA**)

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

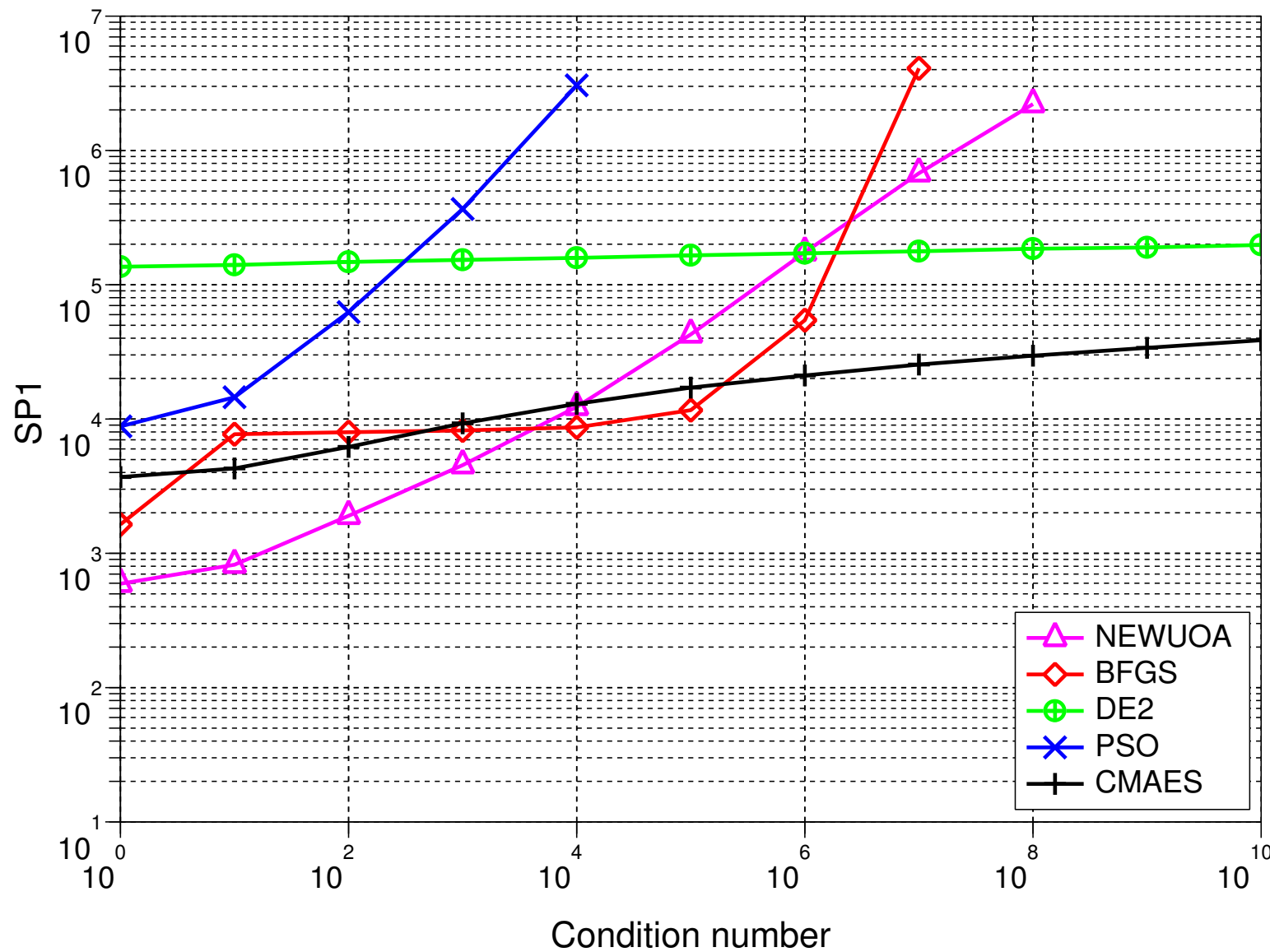
SP1 = average number of objective function evaluations<sup>15</sup> to reach the target function value of  $g^{-1}(10^{-9})$

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

# Comparison to BFGS, NEWUOA, PSO and DE

$f$  non-convex, non-separable (rotated) with varying condition number  $\alpha$

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



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

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

$\mathbf{H}$  full

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

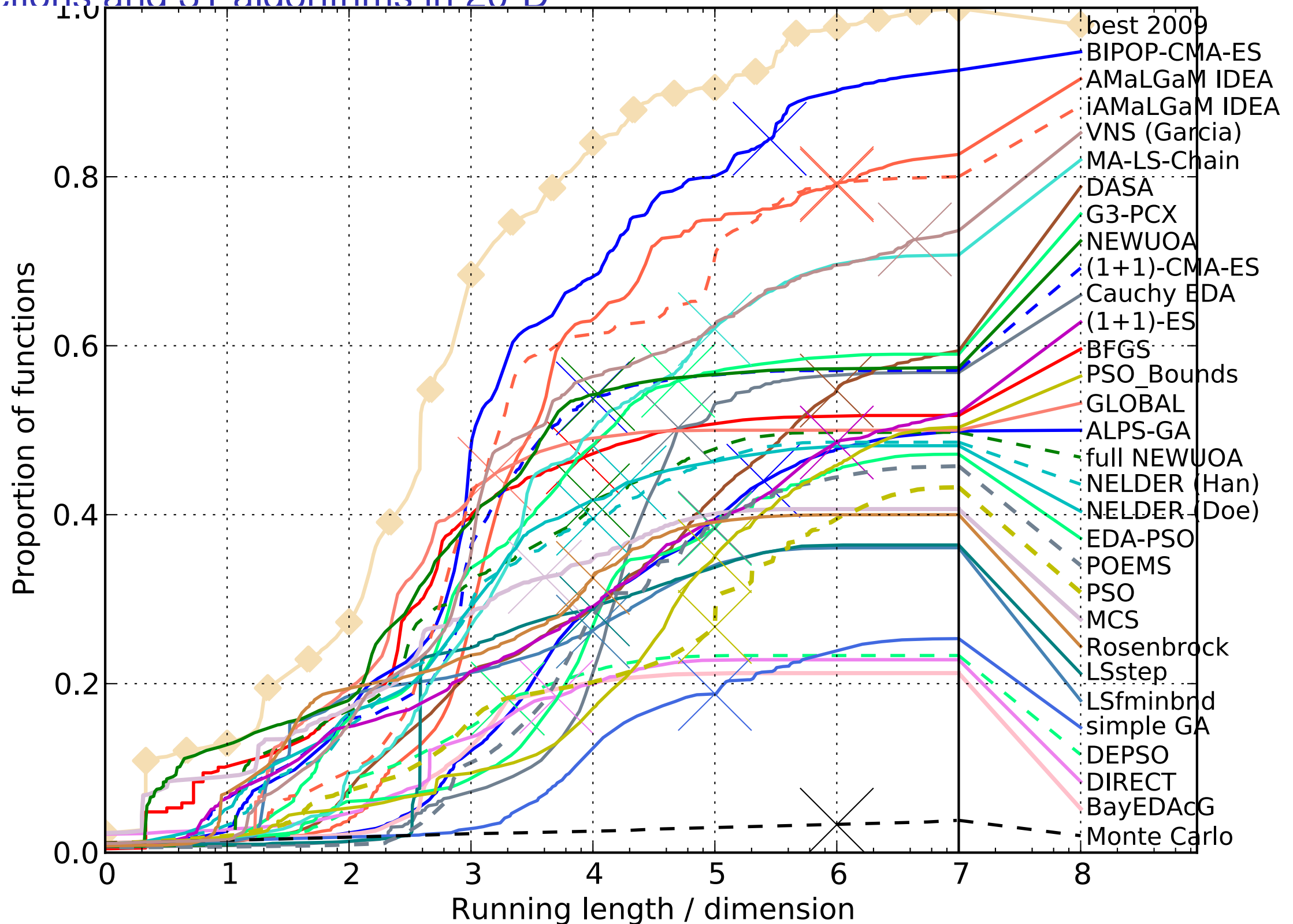
$g$  any order-preserving = strictly increasing function (for all other)

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

<sup>16</sup> Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

# Comparison during BBOB at GECCO 2009

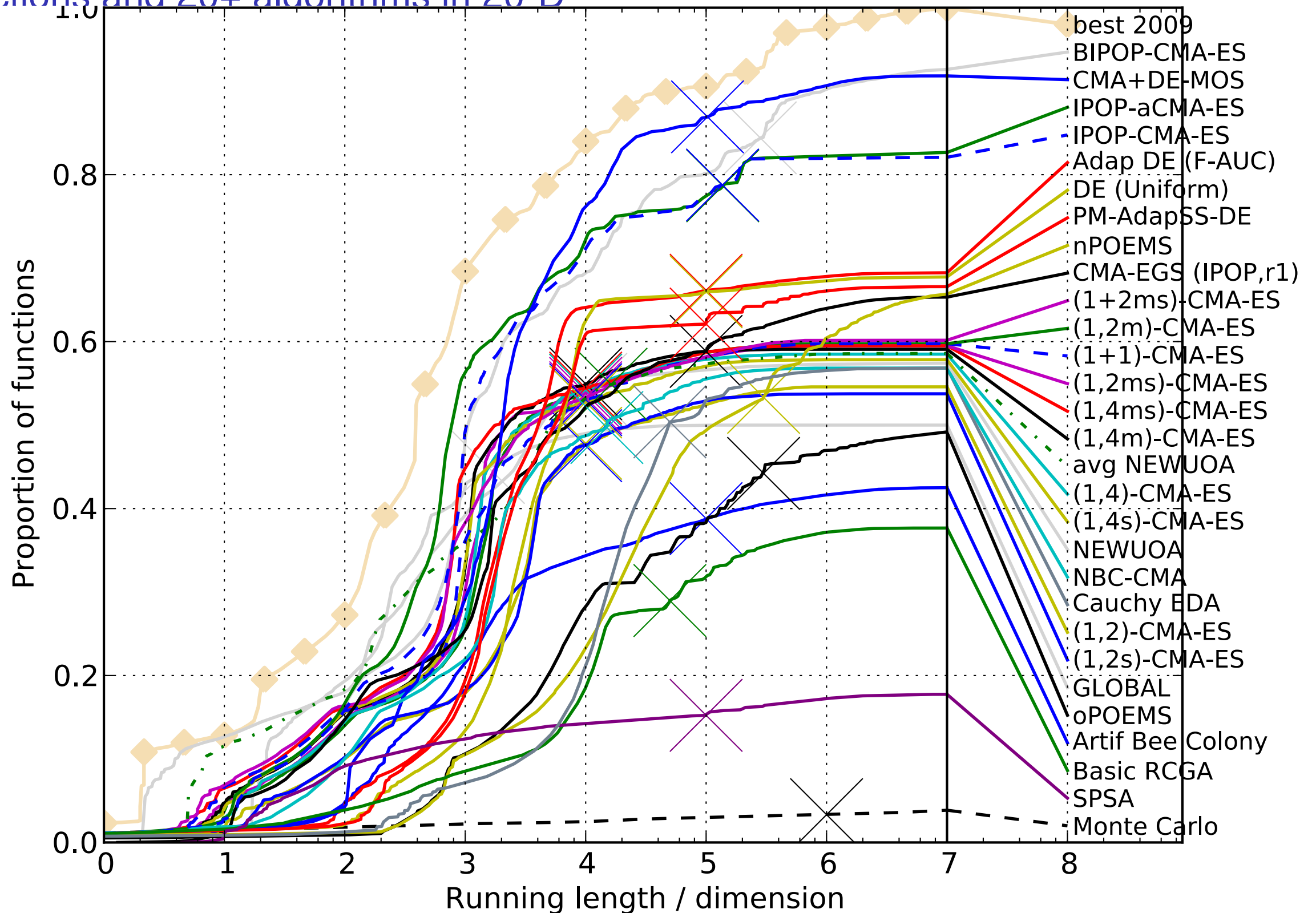
24 functions and 31 algorithms in 20-D





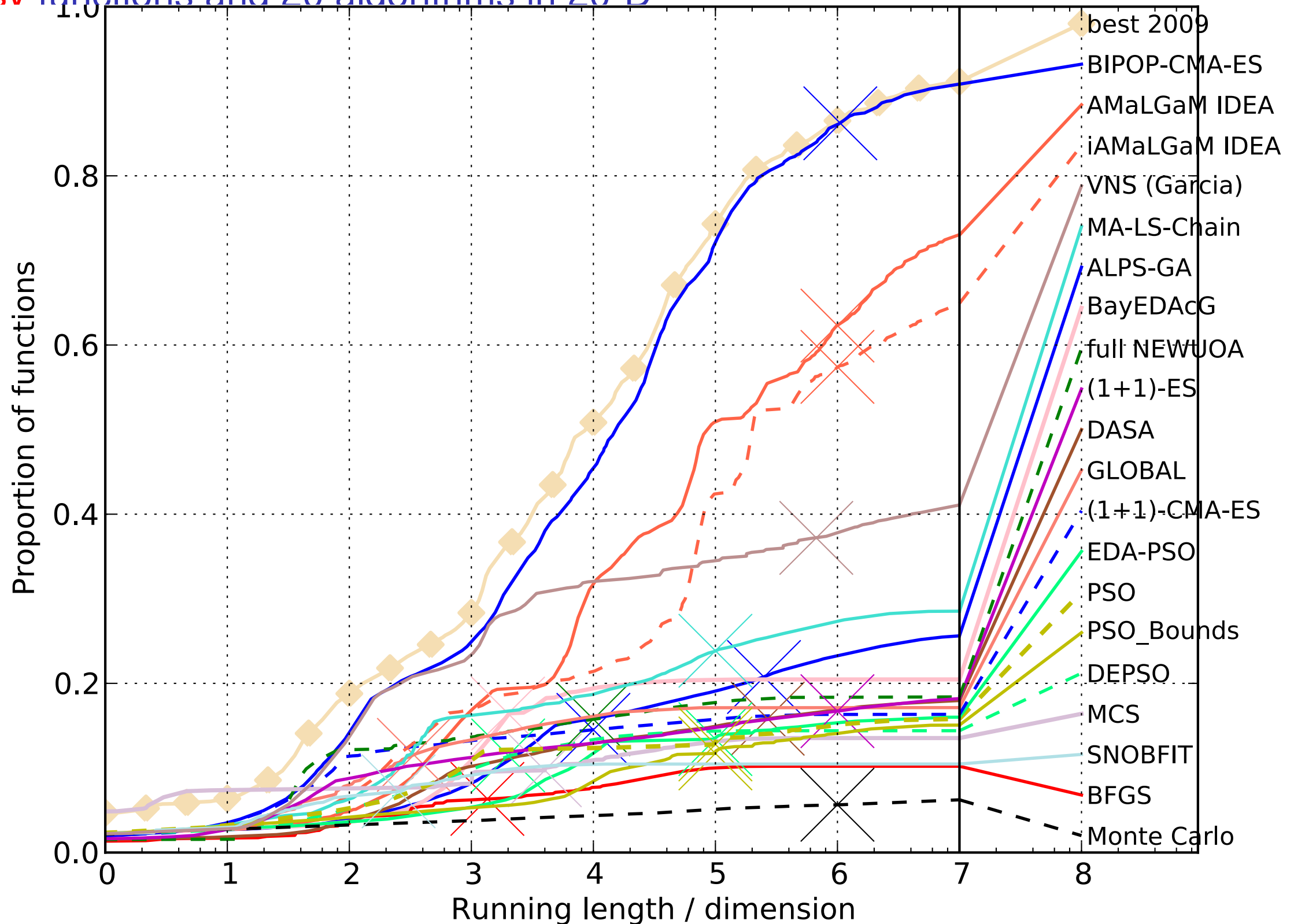
# Comparison during BBOB at GECCO 2010

24 functions and 20+ algorithms in 20-D



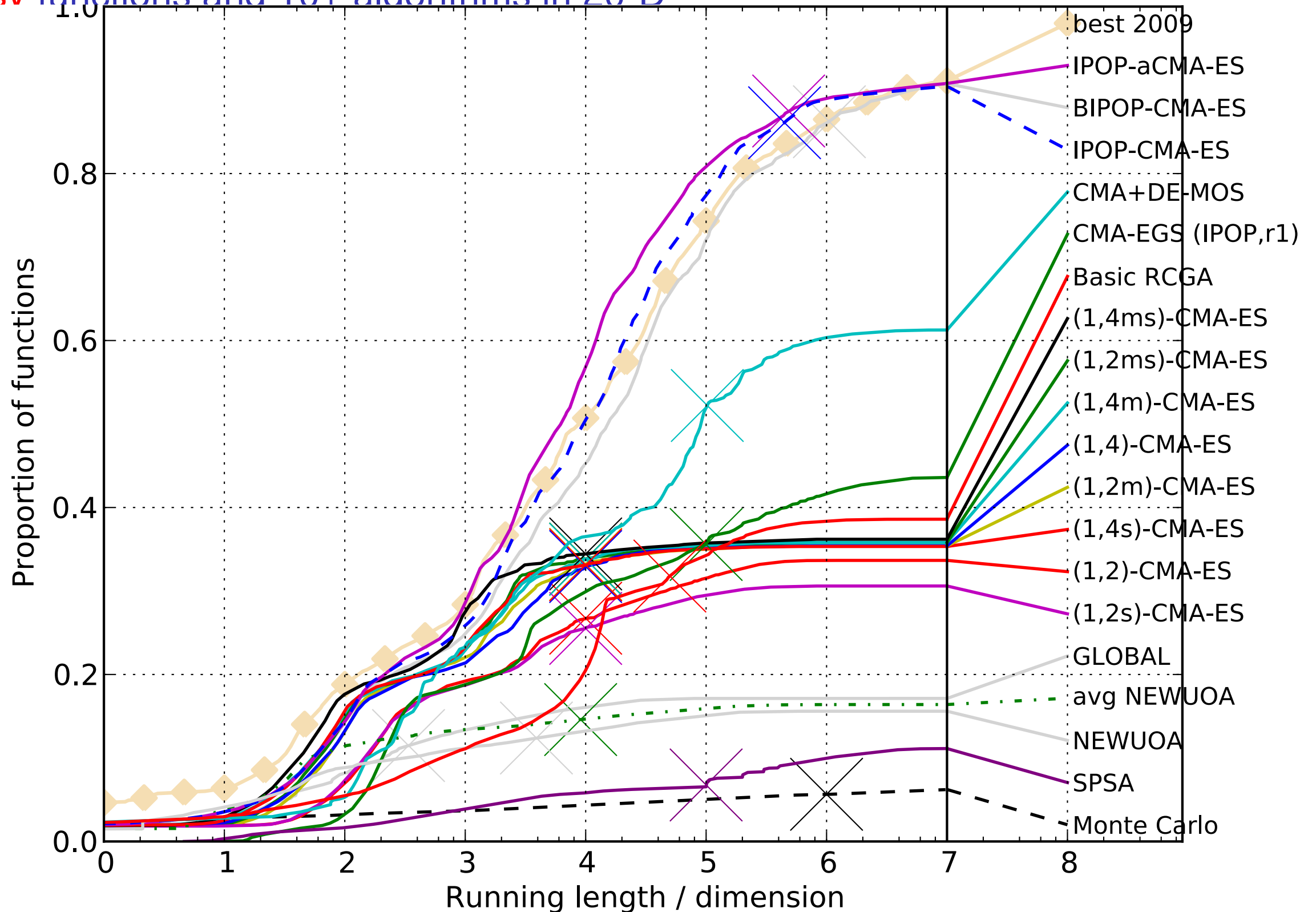
# Comparison during BBOB at GECCO 2009

30 **noisy** functions and 20 algorithms in 20-D



# Comparison during BBOB at GECCO 2010

30 **noisy** functions and 10+ algorithms in 20-D



# Final Remarks

# The Continuous Search Problem

**Difficulties** of a non-linear optimization problem are

- dimensionality and non-separability

demands to exploit problem structure, e.g. neighborhood  
cave: design of benchmark functions

- ill-conditioning

demands to acquire a second order model

- ruggedness

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

# Main Characteristics of (CMA) Evolution Strategies

- 1 Multivariate normal distribution to generate new search points  
follows the maximum entropy principle
- 2 Rank-based selection  
implies invariance, same performance on  $g(f(\mathbf{x}))$  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  
 $\mathbf{C} \propto \mathbf{H}^{-1} \iff$  adapts a variable metric  
 $\iff$  new (rotated) problem representation  
 $\implies f : \mathbf{x} \mapsto g(\mathbf{x}^T \mathbf{H} \mathbf{x})$  reduces to  $\mathbf{x} \mapsto \mathbf{x}^T \mathbf{x}$

# Limitations

## of CMA Evolution Strategies

- **internal CPU-time:**  $10^{-8}n^2$  seconds per function evaluation on a 2GHz PC, tweaks are available  
     1 000 000  $f$ -evaluations in 100-D take 100 seconds *internal* CPU-time
- better methods are presumably available in case of
  - ▶ partly separable problems
  - ▶ specific problems, for example with cheap gradients  
     specific methods
  - ▶ small dimension ( $n \ll 10$ )  
     for example Nelder-Mead
  - ▶ small running times (number of  $f$ -evaluations  $< 100n$ )  
     model-based methods