

Covariance Matrix Adaptation

Covariance Matrix Adaptation

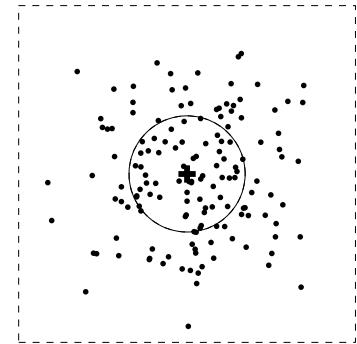
Evolution Strategies

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

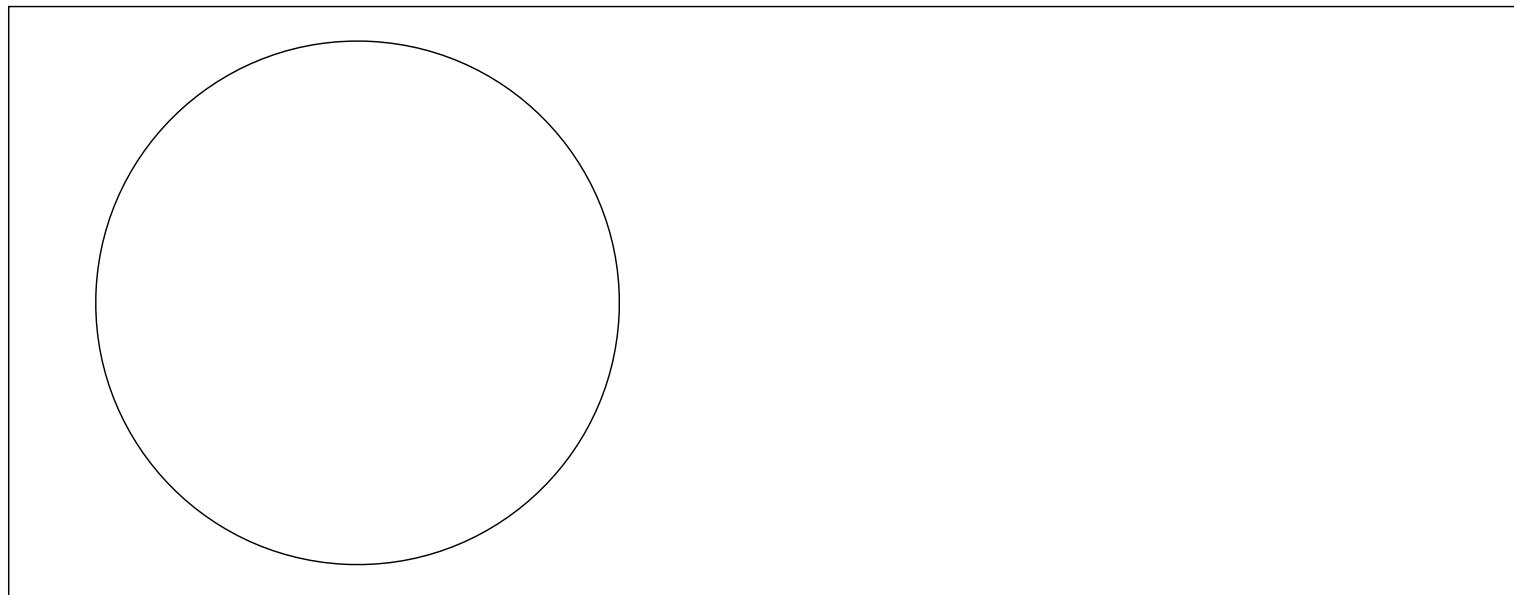
- ▶ the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- ▶ the so-called **step-size** $\sigma \in \mathbb{R}_+$ controls the *step length*
- ▶ the **covariance matrix** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

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

Covariance Matrix Adaptation

Rank-One Update

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

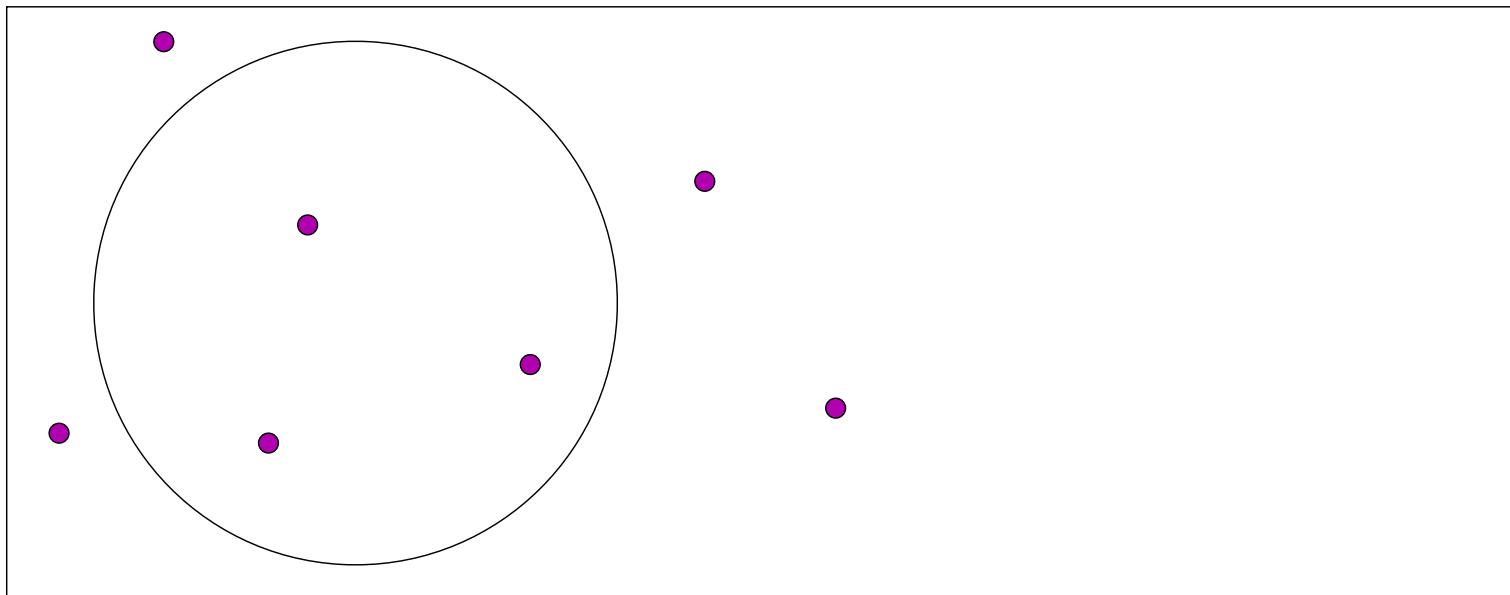


initial distribution, $\mathbf{C} = \mathbf{I}$

Covariance Matrix Adaptation

Rank-One Update

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

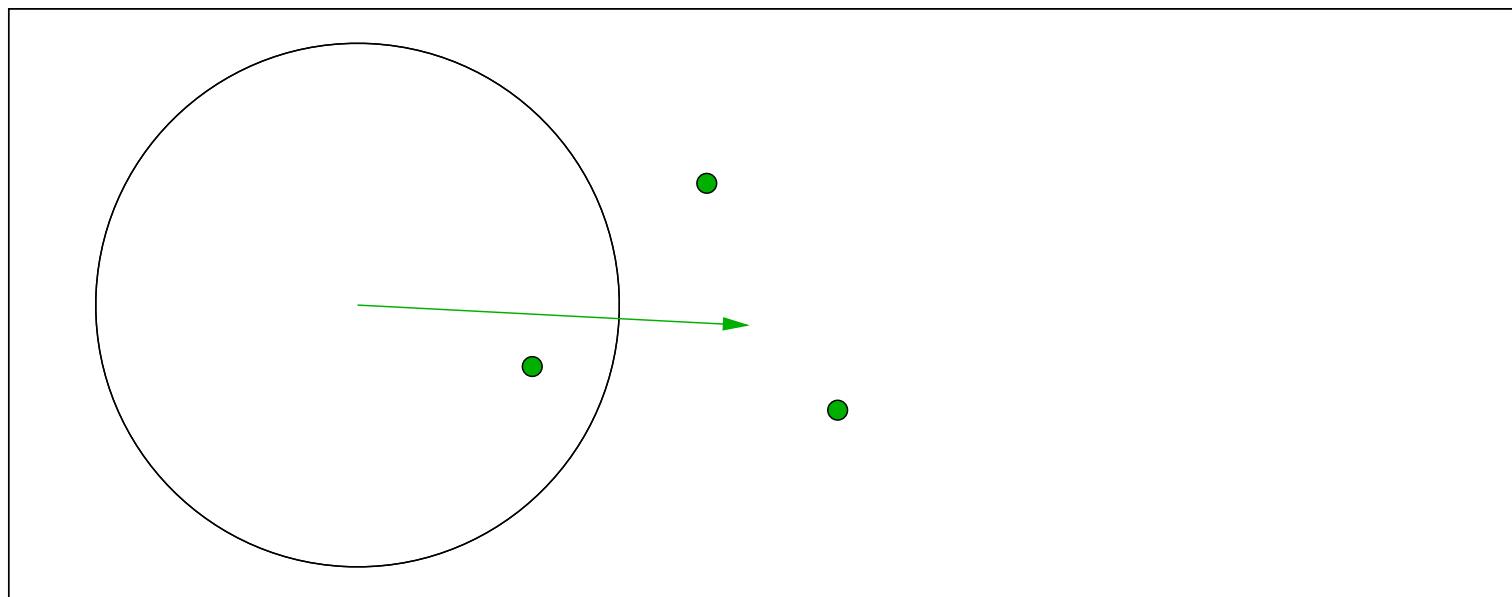


initial distribution, $\mathbf{C} = \mathbf{I}$

Covariance Matrix Adaptation

Rank-One Update

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

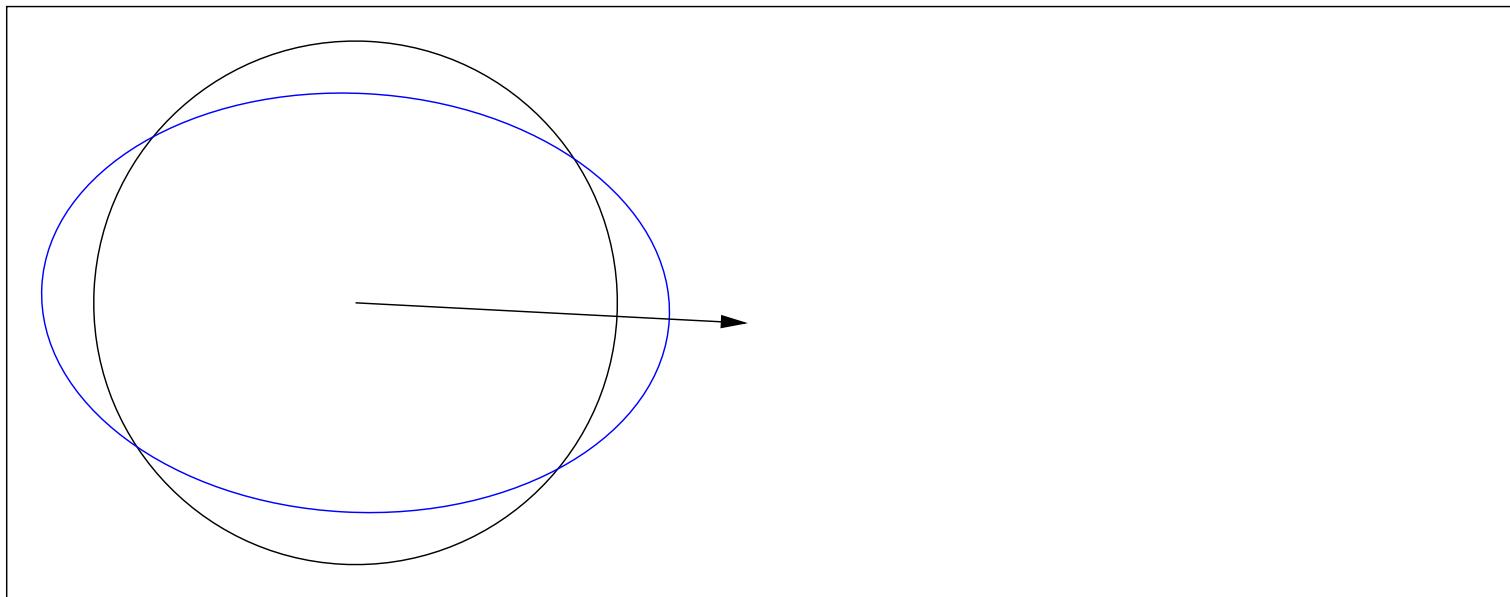


\mathbf{y}_w , movement of the population mean \mathbf{m} (disregarding σ)

Covariance Matrix Adaptation

Rank-One Update

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



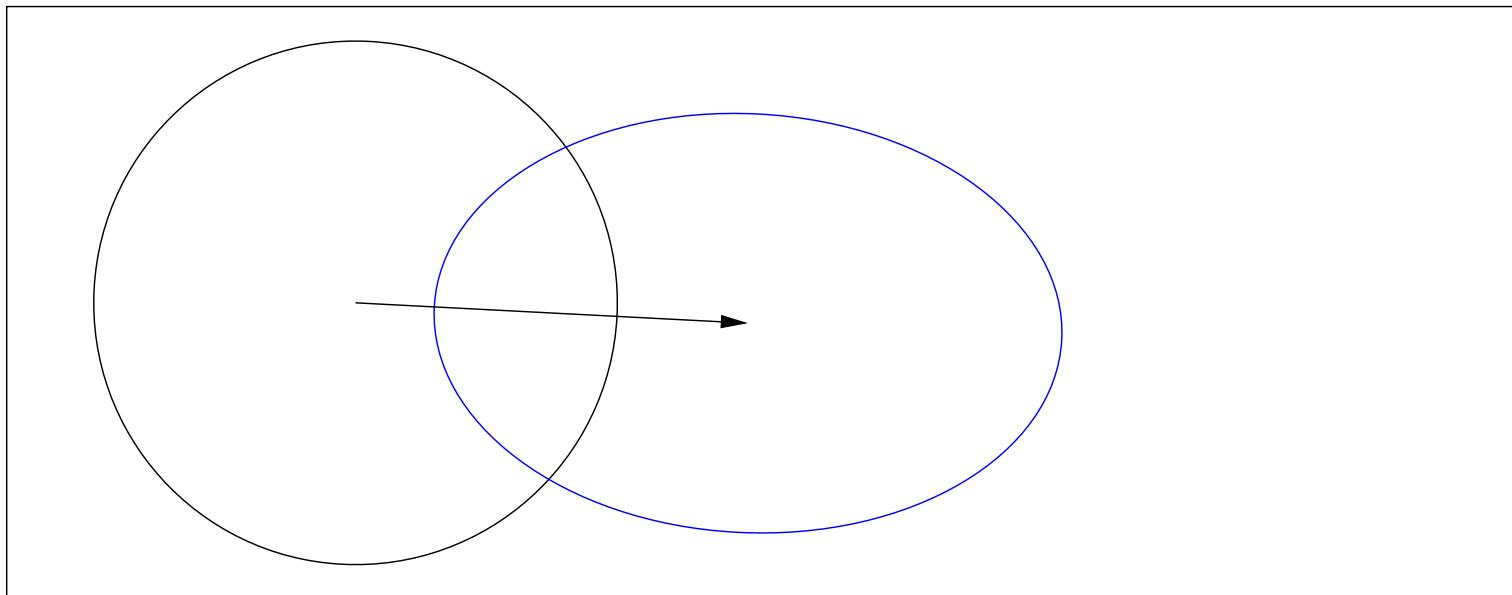
mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

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

Covariance Matrix Adaptation

Rank-One Update

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

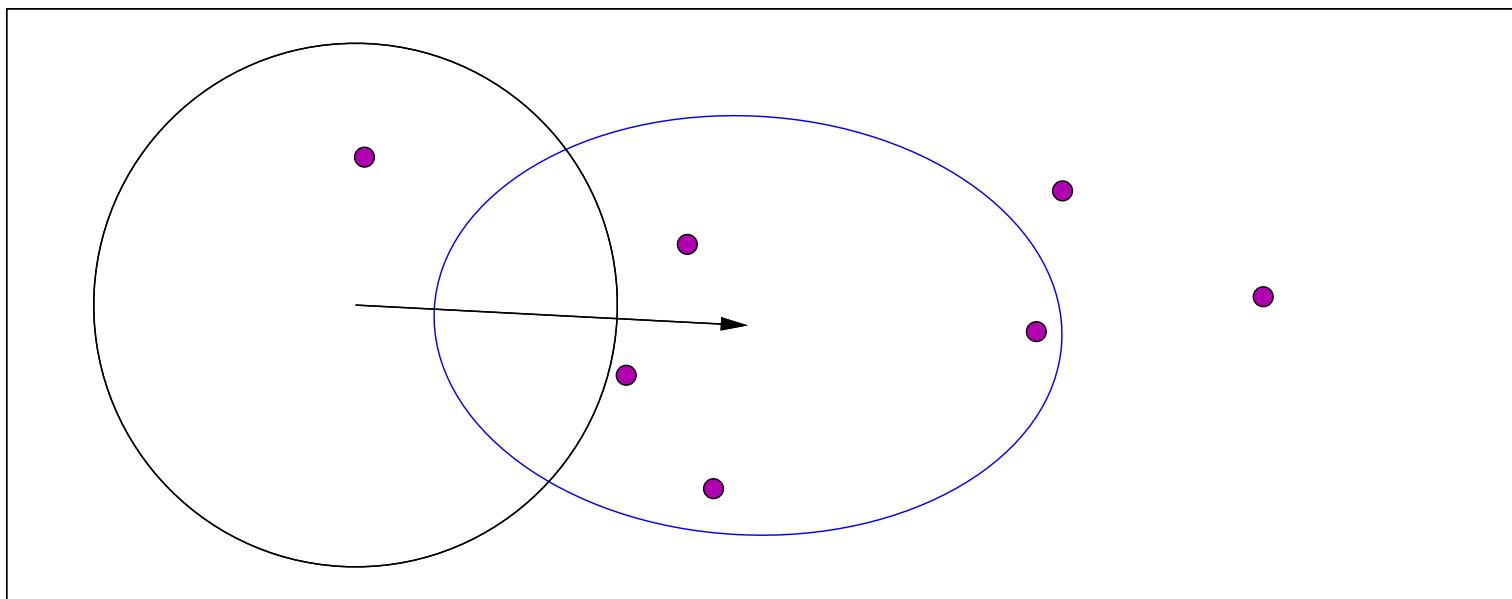


new distribution (disregarding σ)

Covariance Matrix Adaptation

Rank-One Update

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

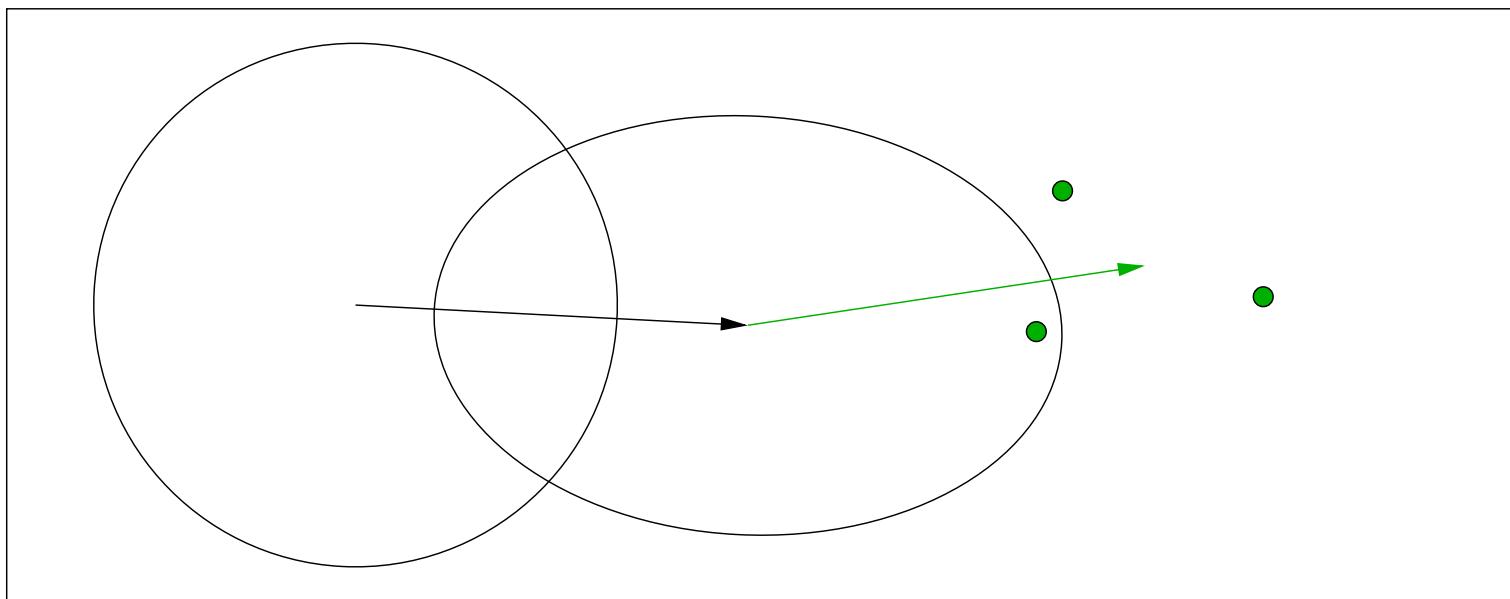


new distribution (disregarding σ)

Covariance Matrix Adaptation

Rank-One Update

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

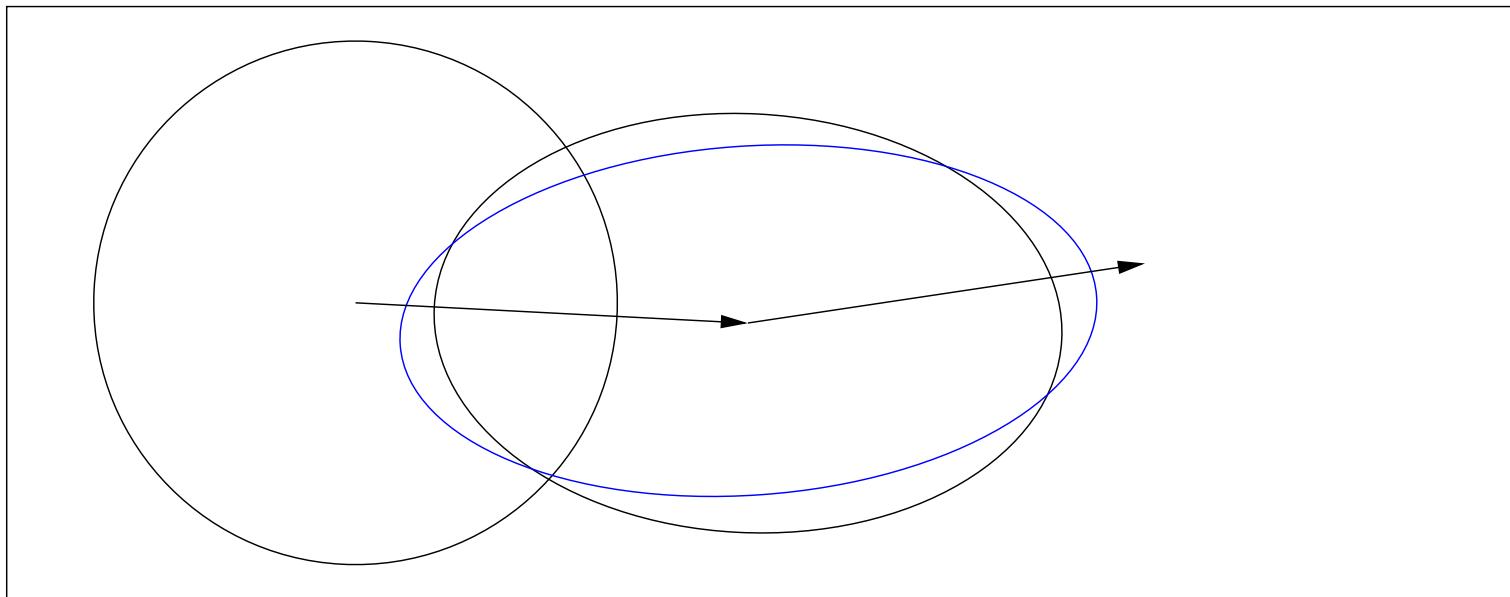


movement of the population mean \boldsymbol{m}

Covariance Matrix Adaptation

Rank-One Update

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



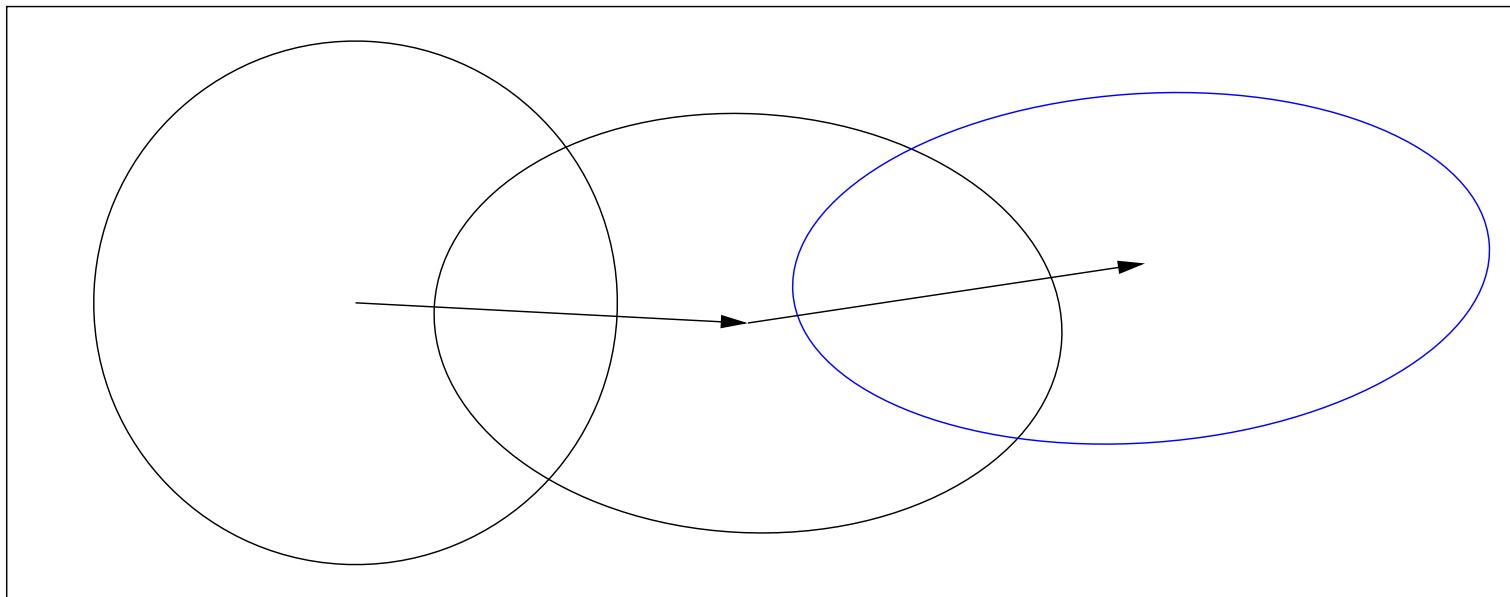
mixture of distribution \mathbf{C} and step \mathbf{y}_w ,

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

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

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

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

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

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

Problem Statement

Stochastic search algorithms - basics

Adaptive Evolution Strategies

Mean Vector Adaptation

Step-size control

Covariance Matrix Adaptation

Rank-One Update

Cumulation—the Evolution Path

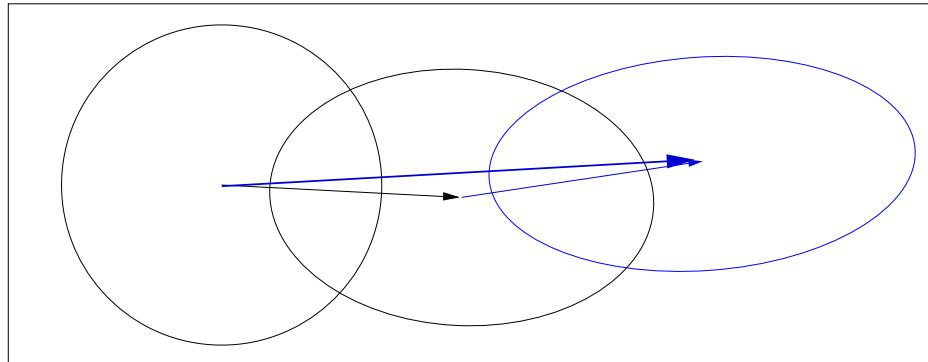
Rank- μ Update

Cumulation

The Evolution Path

Evolution Path

Conceptually, the evolution path is the **search path** the strategy takes over a **number of iteration steps**. It can be expressed as a sum of consecutive *steps* of the mean \mathbf{m} .



An exponentially weighted sum of steps \mathbf{y}_w is used

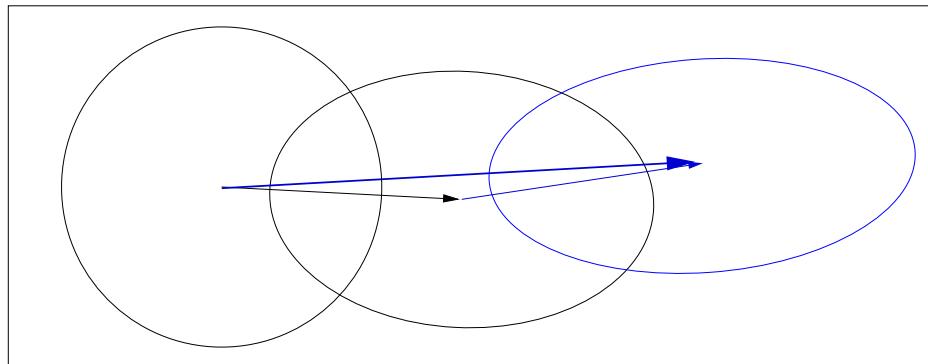
$$\mathbf{p_c} \propto \sum_{i=0}^g \underbrace{(1 - c_c)^{g-i}}_{\text{exponentially fading weights}} \mathbf{y}_w^{(i)}$$

Cumulation

The Evolution Path

Evolution Path

Conceptually, the evolution path is the **search path** the strategy takes over a **number of iteration steps**. It can be expressed as a sum of consecutive *steps* of the mean \mathbf{m} .



An exponentially weighted sum of steps \mathbf{y}_w is used

$$\mathbf{p}_c \propto \sum_{i=0}^g \underbrace{(1 - c_c)^{g-i}}_{\text{exponentially fading weights}} \mathbf{y}_w^{(i)}$$

The recursive construction of the evolution path (cumulation):

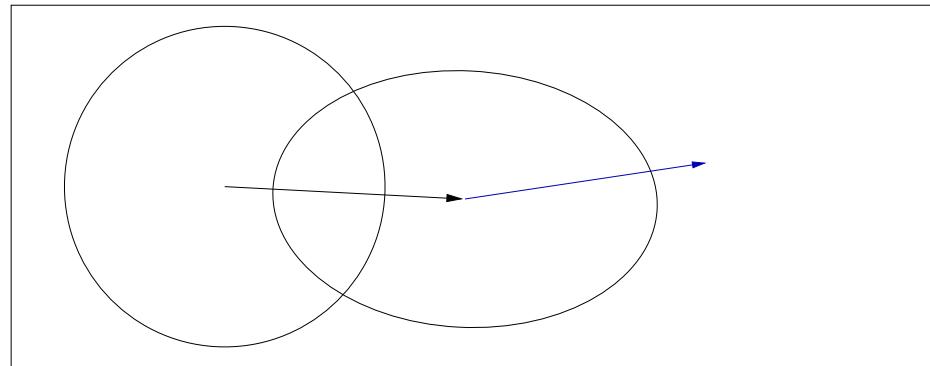
$$\mathbf{p}_c \leftarrow \underbrace{(1 - c_c)}_{\text{decay factor}} \mathbf{p}_c + \underbrace{\sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w}}_{\text{normalization factor}} \underbrace{\mathbf{y}_w}_{\text{input} = \frac{\mathbf{m} - \mathbf{m}_{\text{old}}}{\sigma}}$$

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

Cumulation

Utilizing the Evolution Path

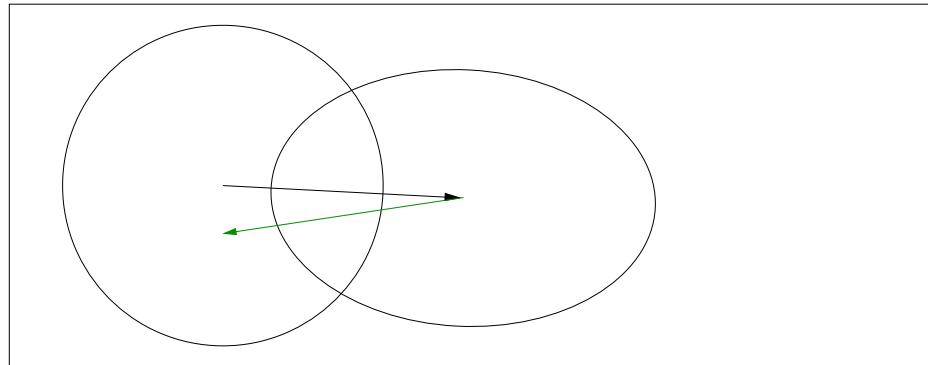
We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



Cumulation

Utilizing the Evolution Path

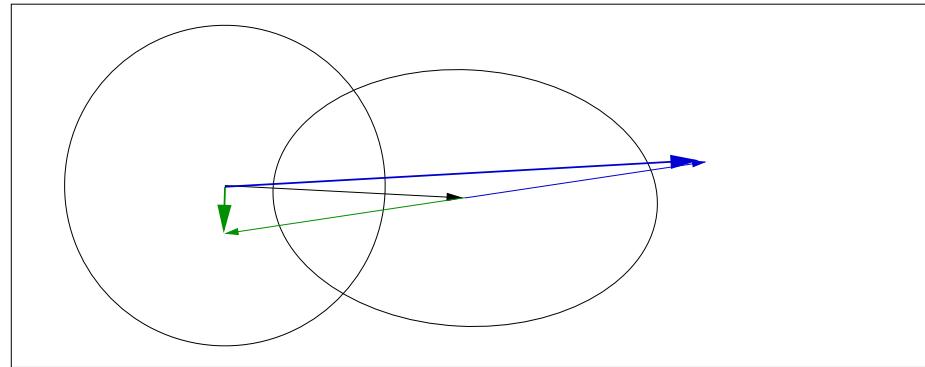
We used $\mathbf{y}_w \mathbf{y}_w^T$ for updating \mathbf{C} . Because $\mathbf{y}_w \mathbf{y}_w^T = -\mathbf{y}_w (-\mathbf{y}_w)^T$ the sign of \mathbf{y}_w is lost.



Cumulation

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

Using an [evolution path](#) for the [rank-one update](#) of the covariance matrix reduces the number of function evaluations to adapt to a straight ridge from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$.⁽³⁾

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

³ Hansen, Müller and Koumoutsakos 2003. Reducing the Time Complexity of the Derandomized Evolution Strategy with Covariance Matrix Adaptation (CMA-ES). *Evolutionary Computation*, 11(1), pp. 1-18

Rank- μ Update

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

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

Rank- μ Update

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

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

The matrix

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

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

Rank- μ Update

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

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

The matrix

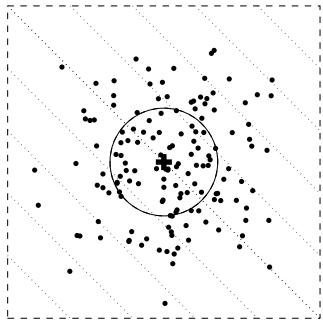
$$\mathbf{C}_\mu = \sum_{i=1}^{\mu} \mathbf{w}_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$$

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

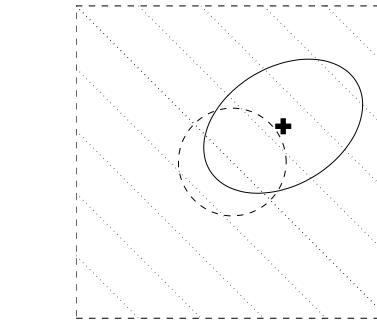
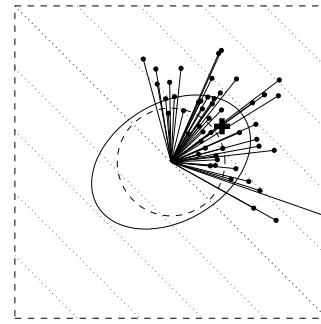
The rank- μ update then reads

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

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



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



$$\mathbf{m}_{\text{new}} \leftarrow \mathbf{m} + \frac{1}{\mu} \sum_{i:\lambda}$$

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
 $c_{\text{cov}} = 1$

new distribution

The rank- μ update

- ▶ increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- ▶ can reduce the number of necessary iterations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁴
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

⁴ 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- μ update

- ▶ increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- ▶ can reduce the number of necessary iterations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁴
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

The rank-one update

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

⁴ 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- μ update

- ▶ increases the possible learning rate in large populations
roughly from $2/n^2$ to μ_w/n^2
- ▶ can reduce the number of necessary iterations roughly from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$ ⁴
given $\mu_w \propto \lambda \propto n$

Therefore the rank- μ update is the primary mechanism whenever a large population size is used

say $\lambda \geq 3n + 10$

The rank-one update

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

Rank-one update and rank- μ update can be combined

⁴ 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: $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, λ

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

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

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

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

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

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

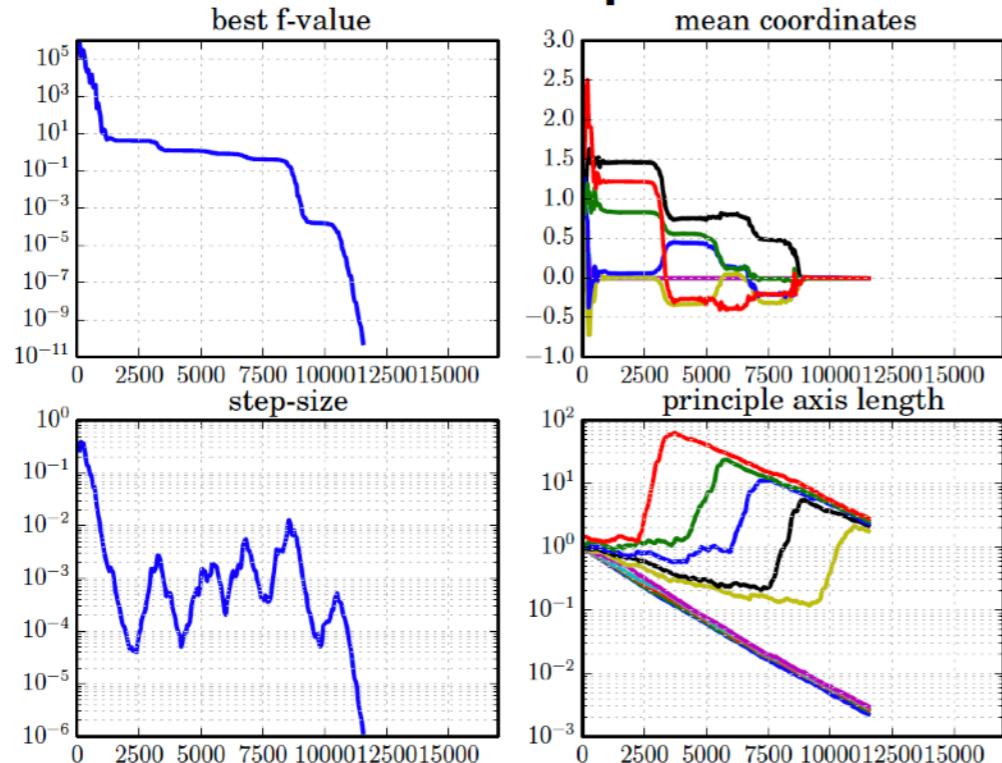
$$\sigma \leftarrow \sigma \times \exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E} \|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right) \quad \text{update of } \sigma$$

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

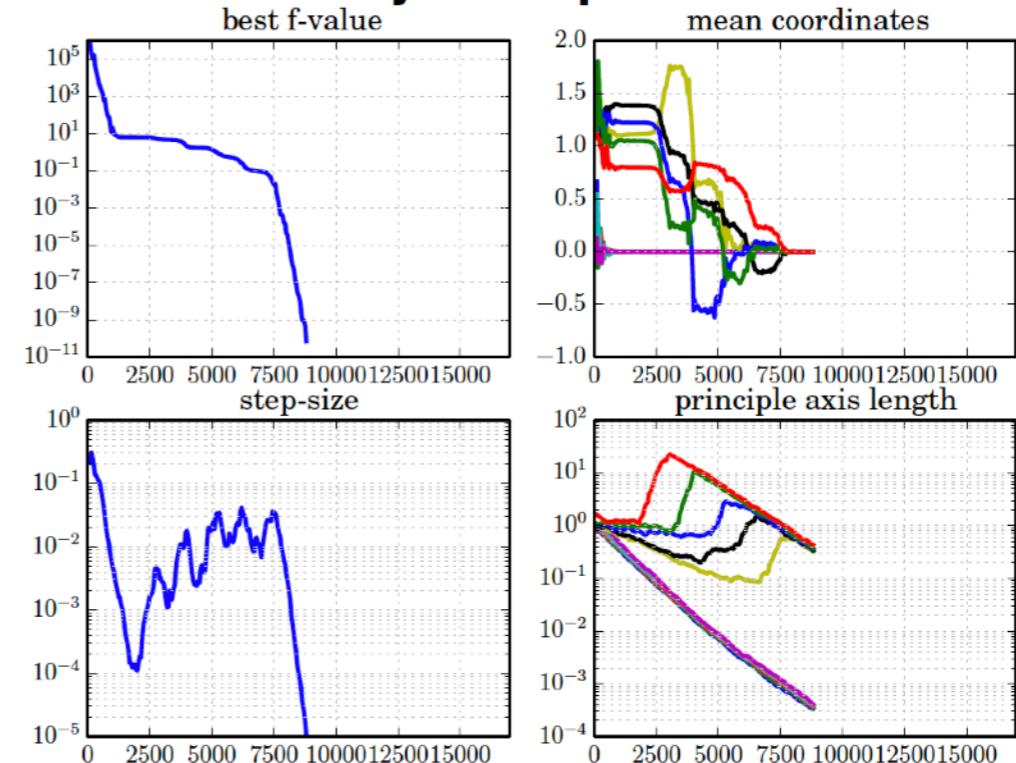
Rank-one and Rank-mu updates

Rank-one and Rank-mu update - default pop size

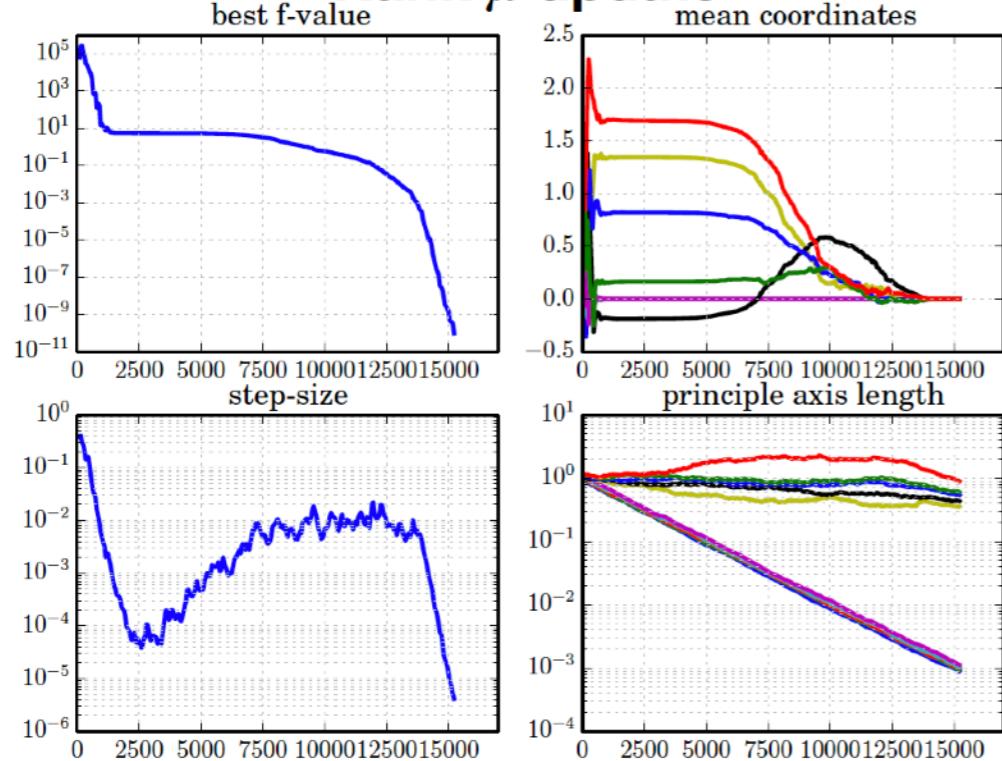
Rank-one update



Hybrid update



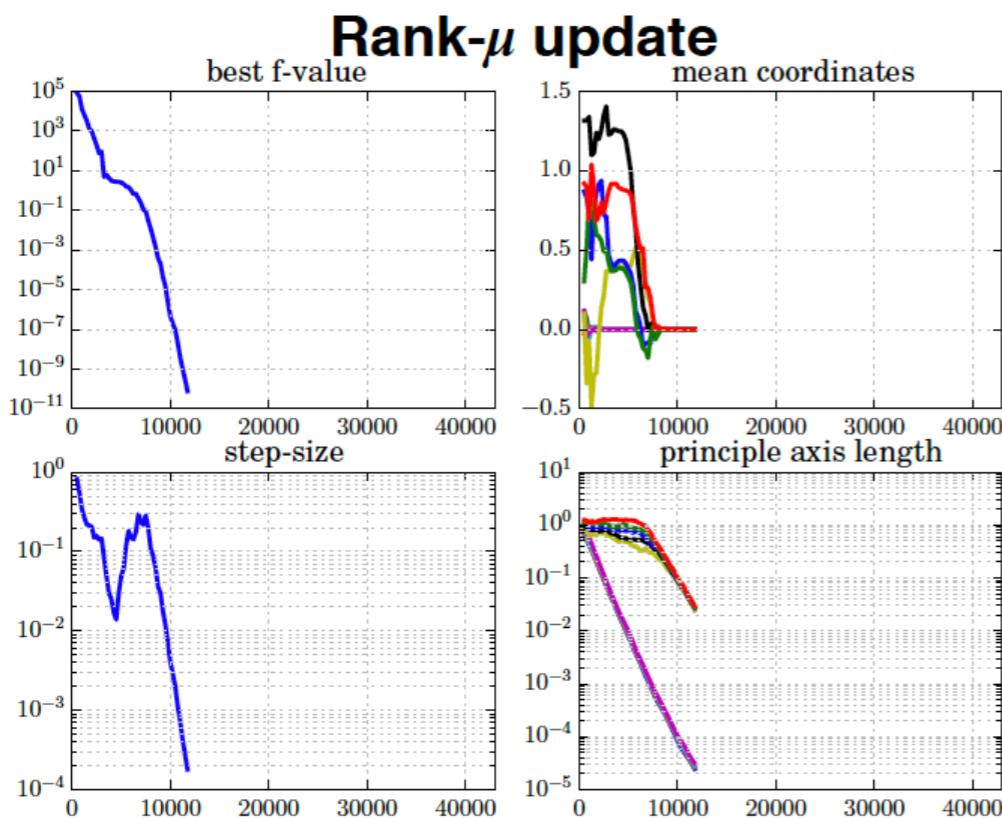
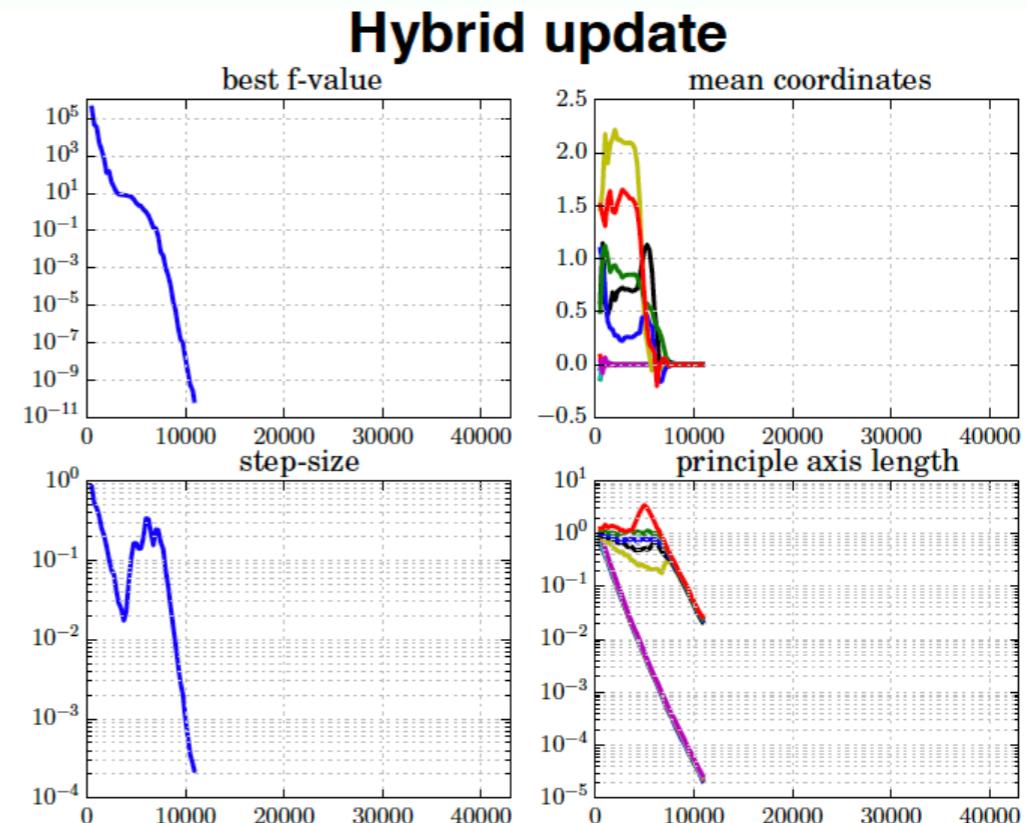
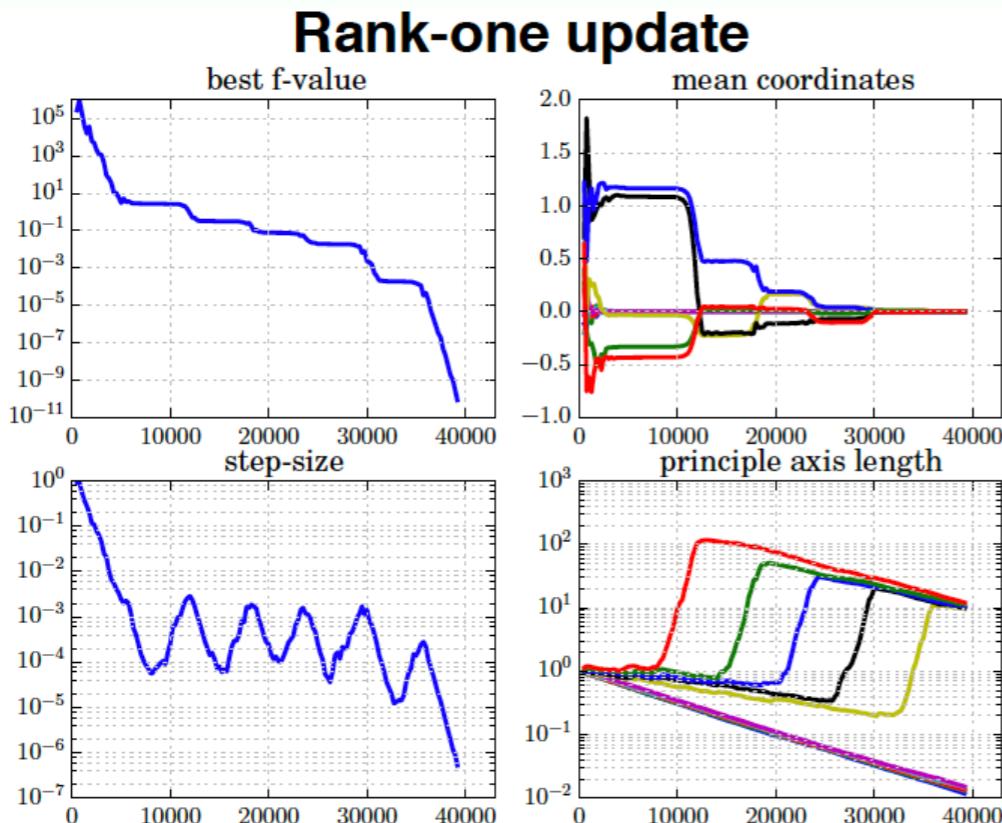
Rank- μ update



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

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

Rank-one and Rank-mu update - larger pop size



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

$$\lambda = 50$$

Experimentum Crucis (0)

What did we want to achieve?

- ▶ reduce any convex-quadratic function

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

to the sphere model

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

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

without use of derivatives

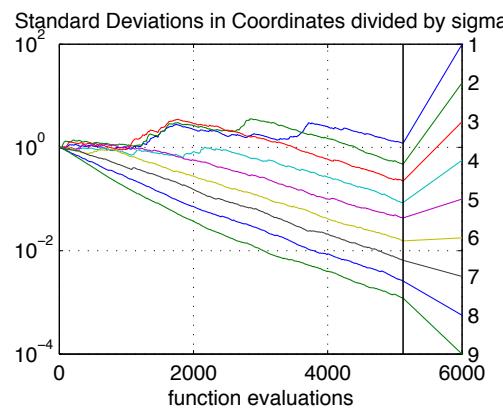
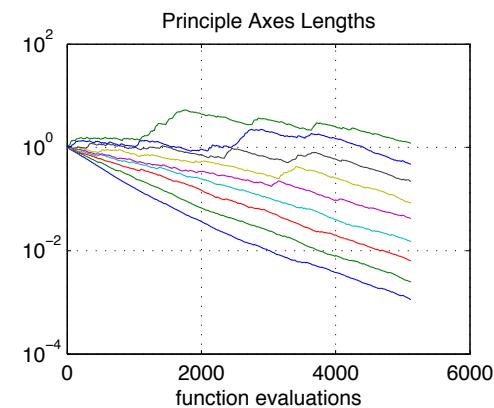
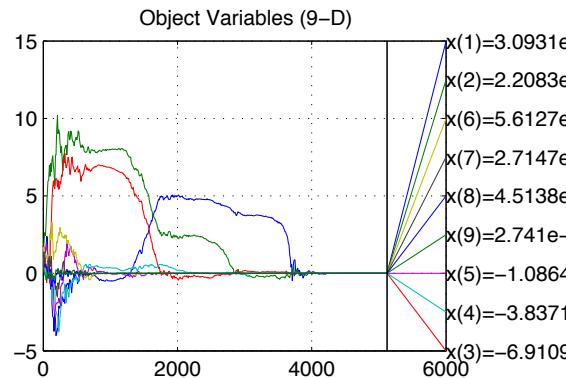
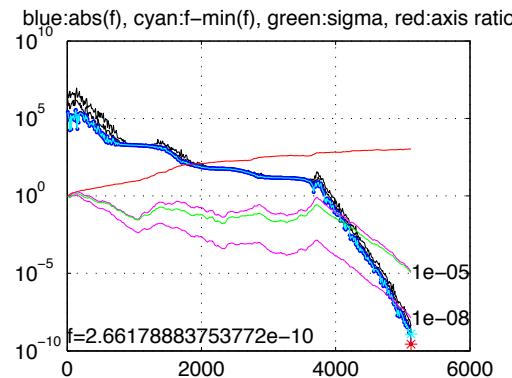
- ▶ lines of equal density align with lines of equal fitness

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

in a stochastic sense

Experimentum Crucis (1)

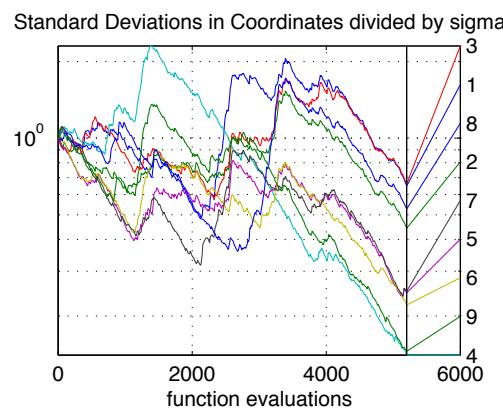
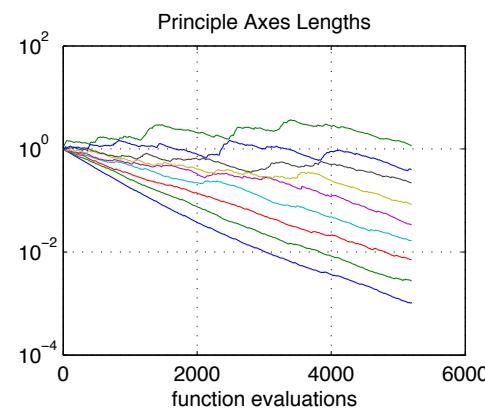
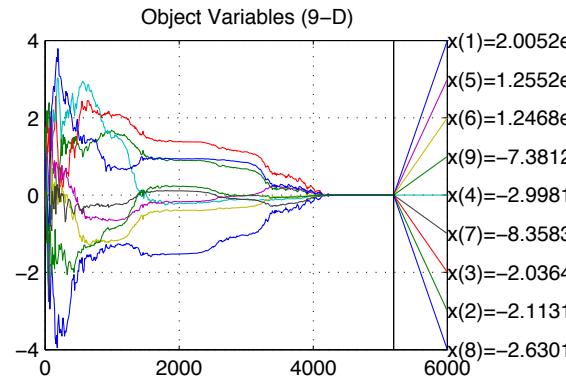
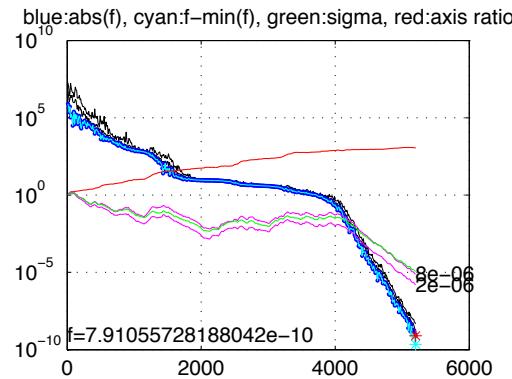
f convex quadratic, separable



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

Experimentum Crucis (2)

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



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

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

On Invariances

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

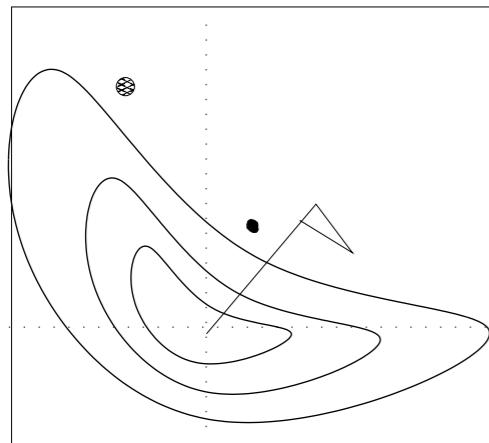
- **Invariance** is a strong **non-empirical** statement about generalization

generalizing (identical) performance from a single function to a whole class of functions

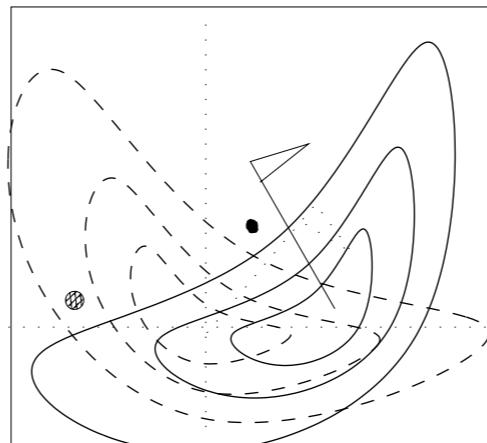
consequently, invariance is important for the evaluation of search algorithms

Rotational Invariance in Search Space

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



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



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

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

45

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

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

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

Main Invariances in Optimization

Invariance to strictly increasing transformations of f : identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto g(f(x)) \text{ where } g : \text{Im}(f) \rightarrow \mathbb{R} \text{ is strictly increasing}$$

Translation invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(x - a) \text{ for all } a \in \mathbb{R}^n$$

Scale invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(\alpha x) \text{ for all } \alpha \in \mathbb{R}_>$$

Rotational invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(Rx) \text{ for all } R \text{ is an orthogonal matrix}$$

Affine invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(Ax + b) \text{ for all } A \in \mathbb{R}^{n \times n} \text{ an invertible matrix and } b \in \mathbb{R}^n$$

Main Invariances in Optimization

Invariance to strictly increasing transformations of f : identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto g(f(x)) \text{ where } g : \text{Im}(f) \rightarrow \mathbb{R} \text{ is strictly increasing}$$

Translation invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(x - a) \text{ for all } a \in \mathbb{R}^n$$

Scale invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(\alpha x) \text{ for all } \alpha \in \mathbb{R}_>$$

provided initial state is change accordingly

Rotational invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(Rx) \text{ for all } R \text{ an orthogonal matrix}$$

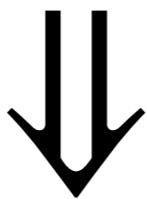
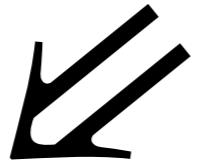
Affine invariance: identical behavior when optimizing

$$x \mapsto f(x)$$

$$x \mapsto f(Ax + b) \text{ for all } A \in \mathbb{R}^{n \times n} \text{ an invertible matrix and } b \in \mathbb{R}^n$$

Hierarchy of Invariance

Affine invariance



Rotational Invariance

Scale-invariance

translation invariance

Exercice - Invariances of (1+1)-ES and CMA-ES

CMA-ES

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

translation invariance

scale invariance

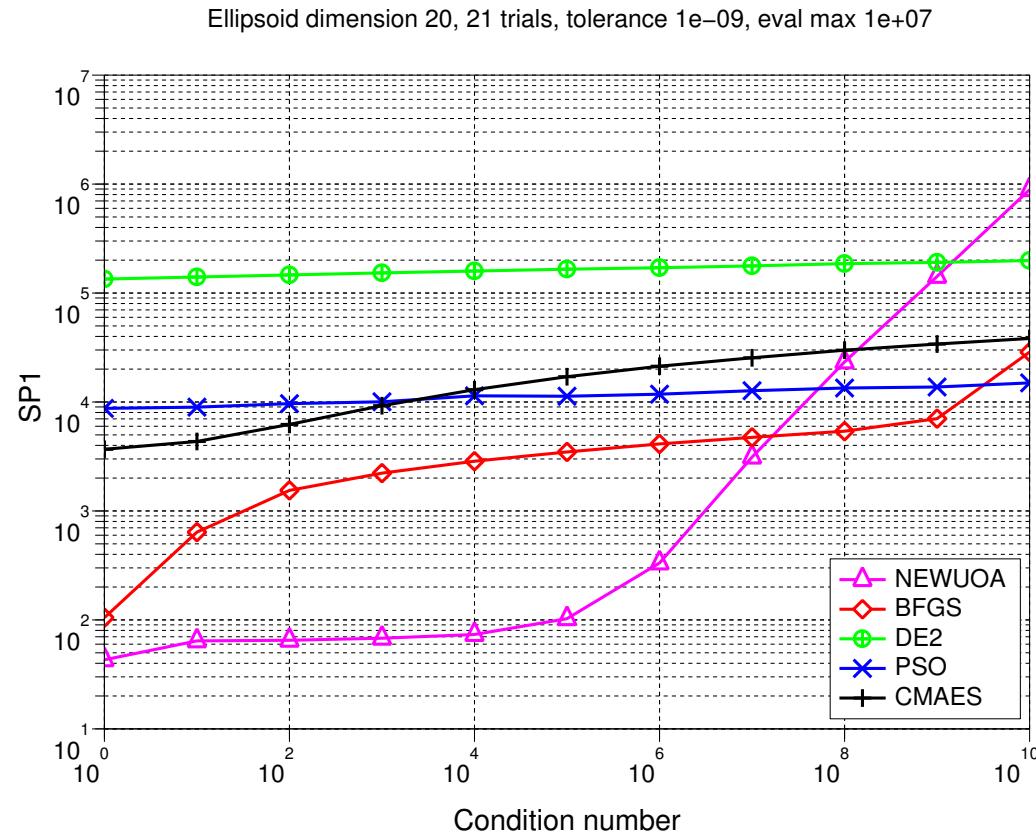
rotational invariance

affine invariance

Testing for invariances

Comparison to BFGS, NEWUOA, PSO and DE

f convex quadratic, separable with varying condition number α



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

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

\mathbf{H} diagonal

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

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

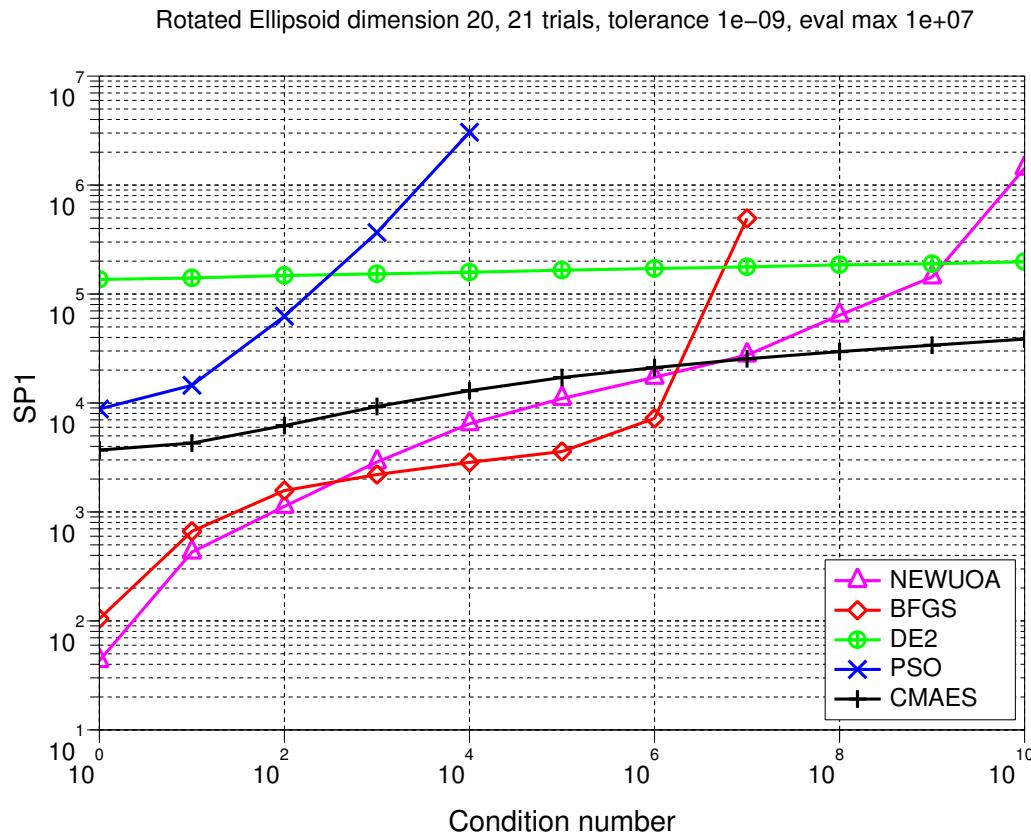
SP1 = average number of objective function evaluations⁵ to reach the target function value of $g^{-1}(10^{-9})$

⁵

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 α



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

$f(x) = g(x^T \mathbf{H} x)$ 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 evaluations⁶ to reach the target function value of $g^{-1}(10^{-9})$

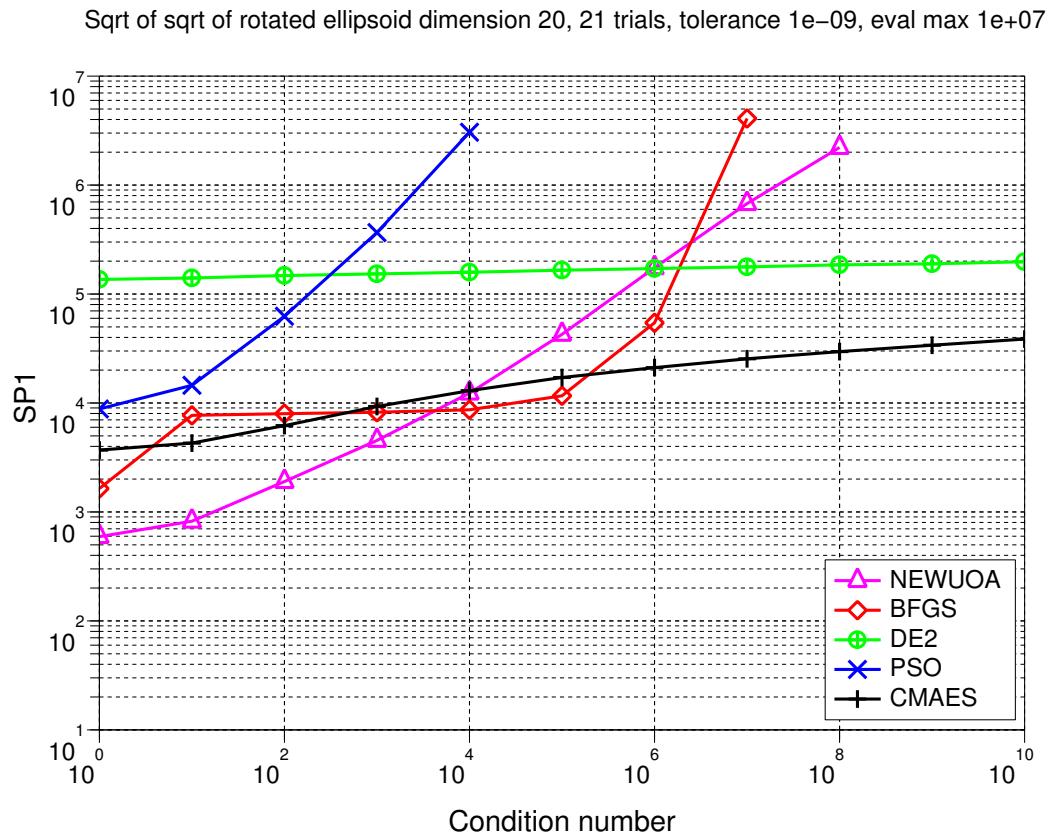
⁶

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 α



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

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

\mathbf{H} full

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

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

SP1 = average number of objective function evaluations⁷ to reach the target function value of $g^{-1}(10^{-9})$

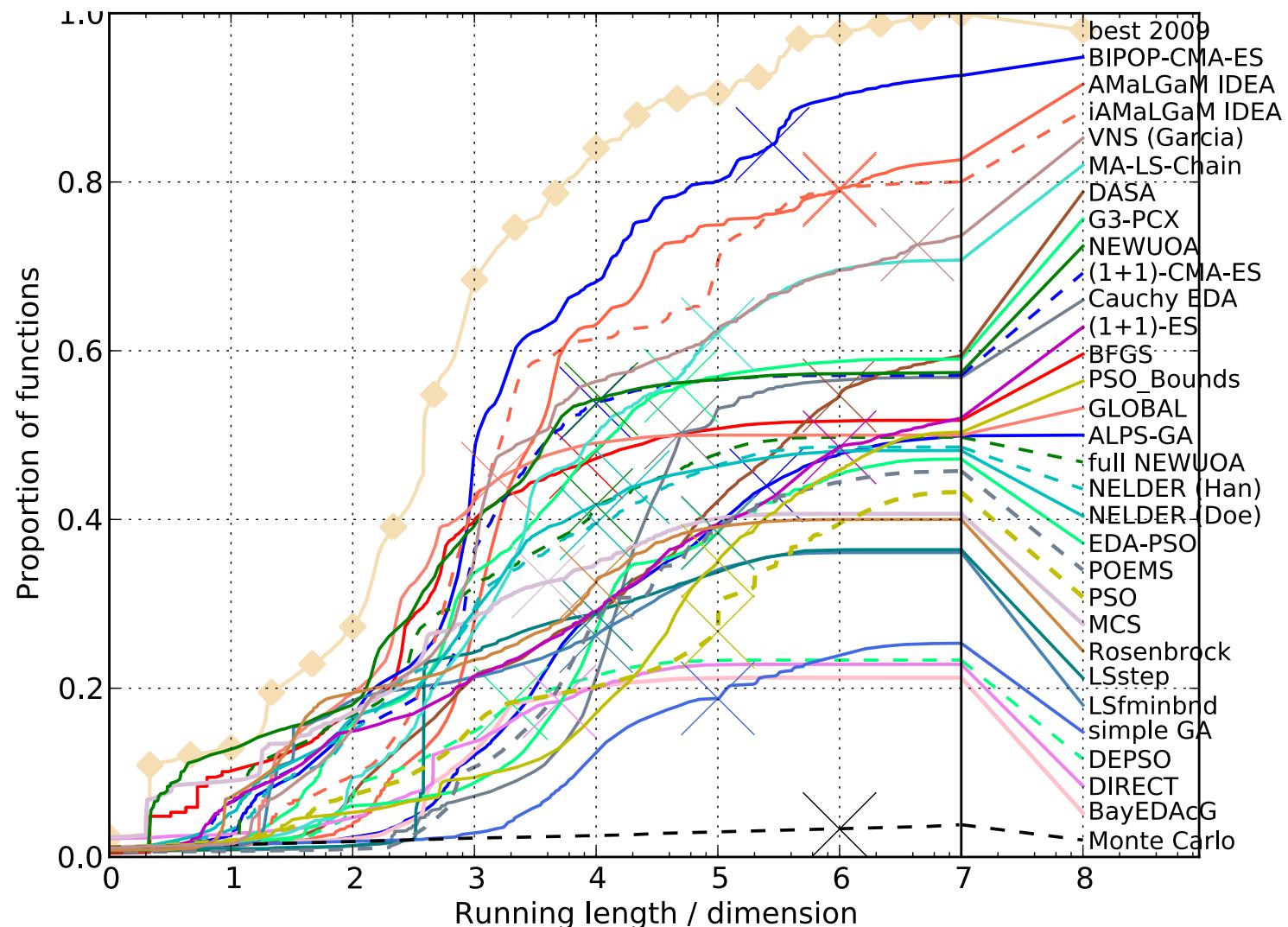
⁷

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



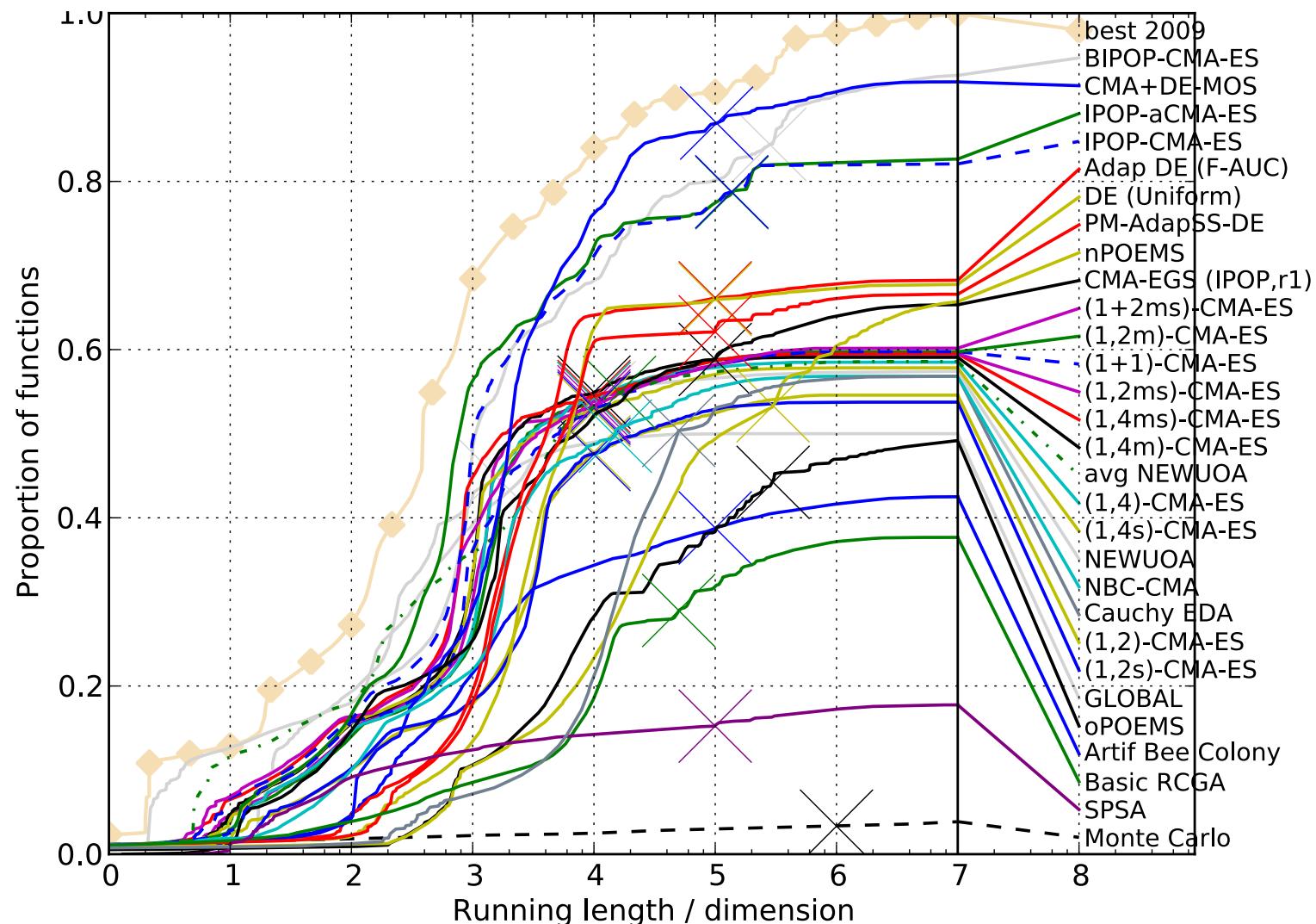
Comparison during BBOB at GECCO 2009

24 functions and 31 algorithms in 20-D



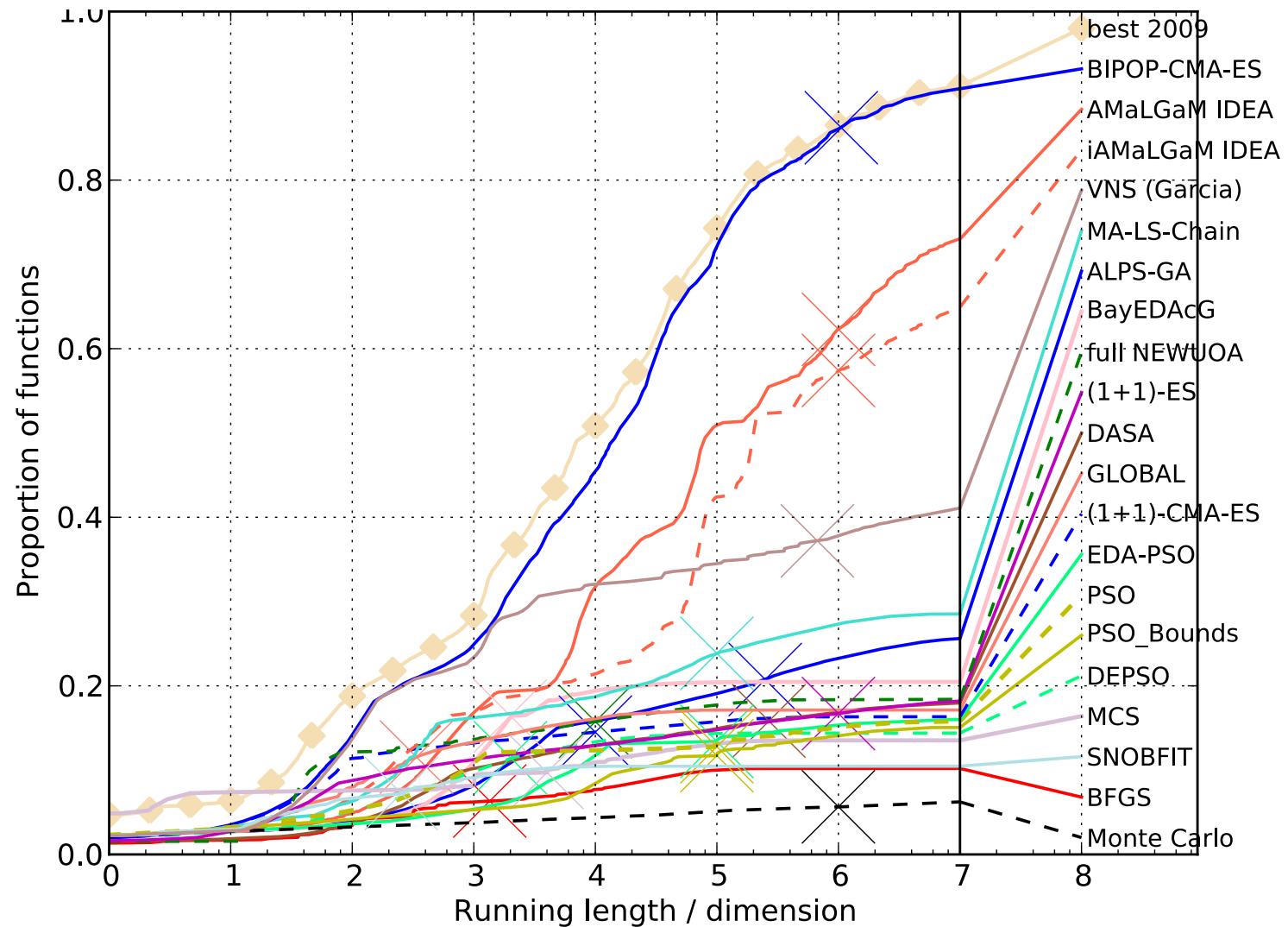
Comparison during BBOB at GECCO 2010

24 functions and 20+ algorithms in 20-D



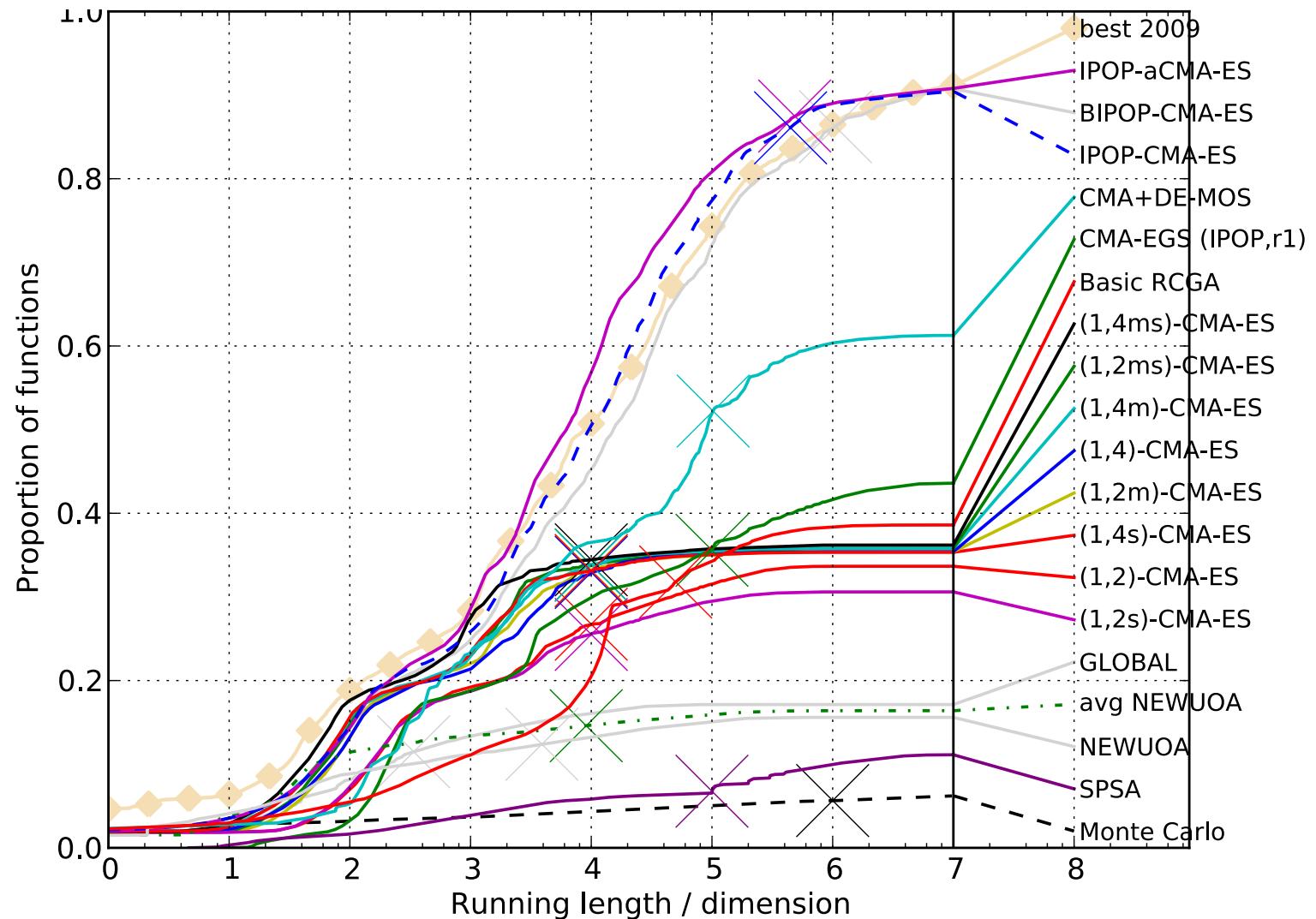
Comparison during BBOB at GECCO 2009

30 noisy functions and 20 algorithms in 20-D



Comparison during BBOB at GECCO 2010

30 noisy functions and 10+ algorithms in 20-D



Problem Statement

Stochastic search algorithms - basics

Adaptive Evolution Strategies

Mean Vector Adaptation

Step-size control

Covariance Matrix Adaptation

Rank-One Update

Cumulation—the Evolution Path

Rank- μ Update

The Continuous Search Problem

Difficulties of a non-linear optimization problem are

- ▶ dimensionality and non-separability
 - demands to exploit problem structure, e.g. neighborhood
- ▶ ill-conditioning
 - demands to acquire a second order model
- ▶ ruggedness
 - demands a non-local (stochastic?) approach

Approach: population based stochastic search, coordinate system independent and with second order estimations (covariances)

Main Features of (CMA) Evolution Strategies

1. Multivariate normal distribution to generate new search points
follows the maximum entropy principle
2. Rank-based selection
implies invariance, same performance on
 $g(f(x))$ for any increasing g
more invariance properties are featured
3. Step-size control facilitates fast (log-linear) convergence
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
 $\mathbf{C} \propto \mathbf{H}^{-1} \iff$ adapts a variable metric
 \iff new (rotated) problem representation
 $\implies f(\mathbf{x}) = g(\mathbf{x}^T \mathbf{H} \mathbf{x})$ reduces to $g(\mathbf{x}^T \mathbf{x})$

Limitations of CMA Evolution Strategies

- ▶ internal CPU-time: $10^{-8} n^2$ seconds per function evaluation on a 2GHz PC, tweaks are available
 - 100 000 f -evaluations in 1000-D take 1/4 hours
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 $\ll 100n$)
model-based methods