

# Algorytm CMA-ES i jego wybrane rozszerzenia

Mateusz Zaborski

M.Zaborski@mini.pw.edu.pl

Faculty of Mathematics and Information Science  
Warsaw University of Technology

16.12.2020

## Table of Content

- 1 Black Box optimization
- 2 Benchmarking
- 3 CMA-ES algorithm
- 4 CMA-ES extensions
  - IPOP-CMA-ES
  - LQ-CMA-ES
- 5 Experiments
  - CMA-ES + Quadratic model + rotation +  $R^2$
  - DE + rotation

# Black Box optimization

## Black Box optimization problem

- Goal - minimize an objective function (or fitness function or cost function)

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

- Single-objective optimization
- Black Box scenario
  - Function values of evaluated search points are the only accessible information
  - Gradients are not available
- Search cost - number of function evaluations

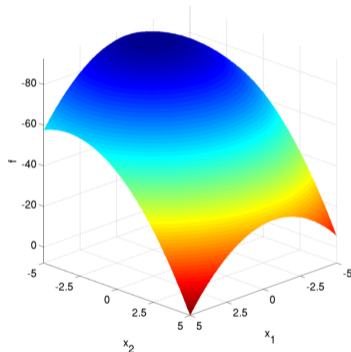


Figure 1: Sphere 2D function from COCO BBOB benchmark [8]

# Black Box optimization

## Applications

### Applications:

- Aerodynamics - shape optimization (simulation-based)
- Economics - e.g. portfolio selection
- Operations research
- Model selection
- Hyperparameter tuning
- etc.

# Black Box optimization

## Performance measures

Performance measures:

- **Number of function evaluations (per target)**
  - Especially for time-consuming evaluations
- Computation time
- Final target reached
- Memory consumption
- Performance for various functions classes

# Black Box optimization

## Problems

Problems:

- Infinite number of solutions in a continuous domain
- Multidimensional problems are difficult for grid search
- The goal is to find a global (not local) optimum
- Unknown function shape
  - Non-linear, non-quadratic
  - Discontinuities, sharp ridges
  - Non-separability
  - Ill-conditioning

# Black Box optimization problem

## Examples

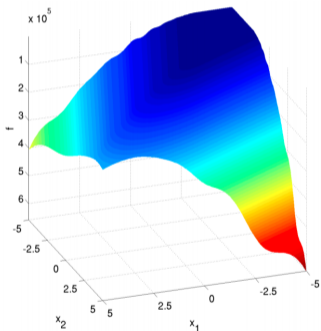


Figure 2: Attractive Sector Function (2D)[8]

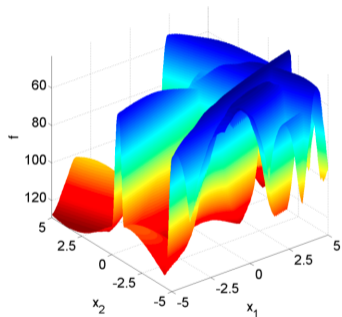


Figure 3: Gallagher's Gaussian 21-hi Peaks Function (2D)[8]



# Black Box optimization problem

## Examples

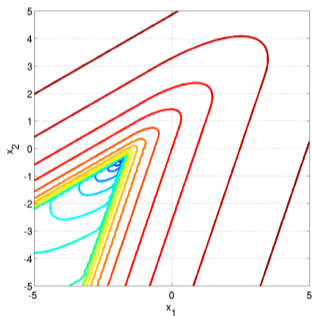


Figure 4: Attractive Sector Function (2D)[8]

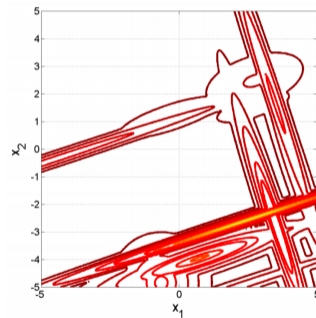


Figure 5: Gallagher's Gaussian 21-hi Peaks Function (2D)[8]

## Optimization techniques

Optimization techniques:

- Random / Monte Carlo
- Grid search
- Simulated annealing
- Bayesian optimization
- Swarm-based algorithms (e.g. PSO)
- Evolutionary algorithms (e.g. DE, CMA-ES)
- Surrogate optimization
- etc.

# Benchmarking

## COCO BBOB benchmark

- Platform for optimizer comparison
- 20–100 functions (24 noiseless all)
- 5–15 instances for each function
- 2, 3, 5, 10, 20, 40 dimensions
- Up to 100 targets per instance

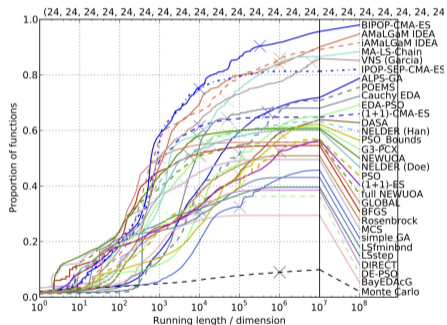


Figure 6: Example empirical runtime distributions from COCO BBOB benchmark [8]

# CMA-ES algorithm

## Covariance Matrix Adaptation Evolution Strategy (CMA-ES)

# Completely Derandomized Self-Adaptation in Evolution Strategies

Nikolaus Hansen and Andreas Ostermeier

In *Evolutionary Computation*, 9(2), pp. 159-195 (2001)

Figure 7: Covariance Matrix Adaptation Evolution Strategy (CMA-ES) idea[9]

# CMA-ES

## Algorithm idea

Initialize distribution parameters  $\theta^{(0)}$

For generation  $g = 0, 1, 2, \dots$

Sample  $\lambda$  independent points from distribution  $P(\mathbf{x}|\theta^{(g)}) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda$

Evaluate the sample  $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$  on  $f$

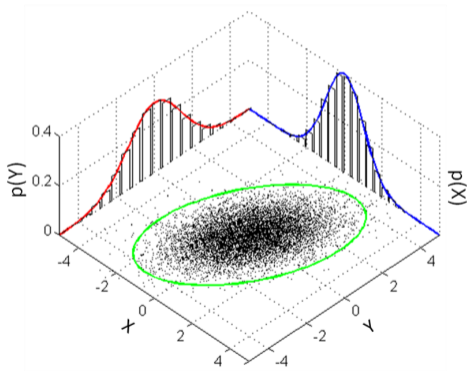
Update parameters  $\theta^{(g+1)} = F_\theta(\theta^{(g)}, (\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_\lambda, f(\mathbf{x}_\lambda)))$

break, if termination criterion met

Figure 8: CMA-ES as randomized black box search algorithm[6]

# CMA-ES

## Multivariate normal distribution



$$m = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

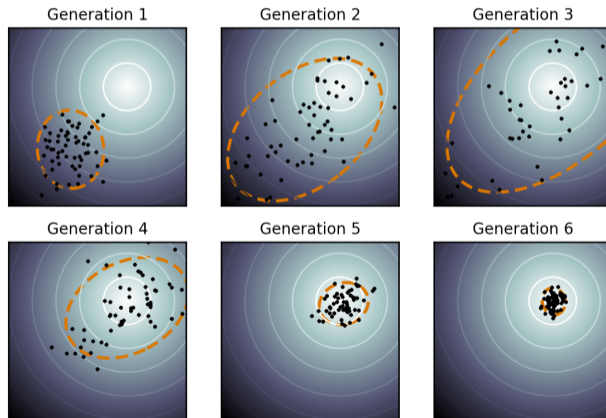
$$C = \begin{bmatrix} 1 & 3/5 \\ 3/5 & 2 \end{bmatrix}$$

Figure 9: Source: [4]



# CMA-ES

## Algorithm idea [3]



# CMA-ES

## Sampling[6]

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

$\mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)})$  is a multivariate normal distribution with zero mean and covariance matrix  $\mathbf{C}^{(g)}$ , see Sect. 0.2. It holds  $\mathbf{m}^{(g)} + \sigma^{(g)} \mathcal{N}(\mathbf{0}, \mathbf{C}^{(g)}) \sim \mathcal{N}(\mathbf{m}^{(g)}, (\sigma^{(g)})^2 \mathbf{C}^{(g)})$ .

$\mathbf{x}_k^{(g+1)} \in \mathbb{R}^n$ ,  $k$ -th offspring (individual, search point) from generation  $g + 1$ .

$\mathbf{m}^{(g)} \in \mathbb{R}^n$ , mean value of the search distribution at generation  $g$ .

$\sigma^{(g)} \in \mathbb{R}_{>0}$ , “overall” standard deviation, step-size, at generation  $g$ .

$\mathbf{C}^{(g)} \in \mathbb{R}^{n \times n}$ , covariance matrix at generation  $g$ . Up to the scalar factor  $\sigma^{(g)2}$ ,  $\mathbf{C}^{(g)}$  is the covariance matrix of the search distribution.

$\lambda \geq 2$ , population size, sample size, number of offspring.

# CMA-ES

## Adaptive components

Adaptive components:

- Mean
- Covariance matrix
  - rank-*one*-update (evolution path  $p_c$ )
  - rank- $\mu$ -update
- Step-size
- Population size (in *IPOP-CMA-ES*[1])

# CMA-ES

## Mean update[6]

$$\mathbf{m}^{(g+1)} = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}^{(g+1)}$$

$$\sum_{i=1}^{\mu} w_i = 1, \quad w_1 \geq w_2 \geq \dots \geq w_{\mu} > 0$$

---

$$\mathbf{m}^{(g+1)} = \mathbf{m}^{(g)} + c_m \sum_{i=1}^{\mu} w_i (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)})$$

Usually  $c_m = 1$

# CMA-ES

## Covariance matrix update

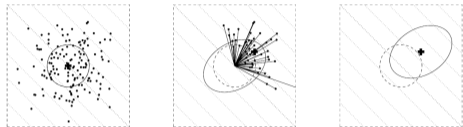


Figure 10: Rank- $\mu$ -update idea[5]

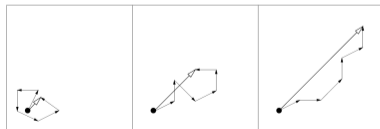


Figure 11: Rank-*one*-update (evolution path) idea[5]

# CMA-ES

## Covariance matrix update[6]

$$\begin{aligned}
 \mathbf{C}^{(g+1)} &= \underbrace{(1 - c_1 - c_\mu \sum w_j)}_{\text{can be close or equal to 0}} \mathbf{C}^{(g)} \\
 &\quad + c_1 \underbrace{\mathbf{p}_c^{(g+1)} \mathbf{p}_c^{(g+1)\top}}_{\text{rank-one update}} + c_\mu \underbrace{\sum_{i=1}^{\lambda} w_i \mathbf{y}_{i:\lambda}^{(g+1)} (\mathbf{y}_{i:\lambda}^{(g+1)})^\top}_{\text{rank-}\mu \text{ update}}
 \end{aligned}$$

$$\mathbf{p}_c^{(g+1)} = (1 - c_c) \mathbf{p}_c^{(g)} + \sqrt{c_c(2 - c_c)\mu_{\text{eff}}} \frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}}$$

$$c_1 \approx 2/n^2$$

$$c_\mu \approx \min(\mu_{\text{eff}}/n^2, 1 - c_1)$$

$$\mathbf{y}_{i:\lambda}^{(g+1)} = (\mathbf{x}_{i:\lambda}^{(g+1)} - \mathbf{m}^{(g)})/\sigma^{(g)}$$

$$\sum w_j = \sum_{i=1}^{\lambda} w_i \approx -c_1/c_\mu$$

# CMA-ES

## Step size

Measure the length of the *evolution path*

the pathway of the mean vector  $m$  in the generation sequence

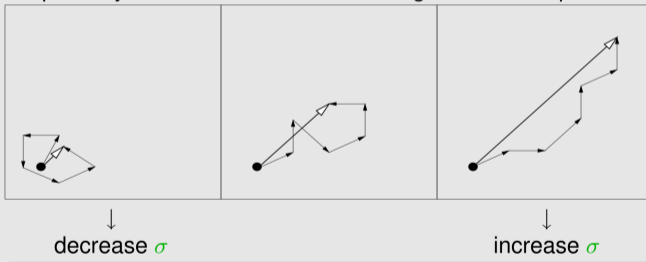


Figure 12: Step size adaptation idea[2]

# CMA-ES

## Conjugate evolution path[6]

$$\mathbf{p}_\sigma^{(g+1)} = (1 - c_\sigma)\mathbf{p}_\sigma^{(g)} + \sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}} \mathbf{C}^{(g)-\frac{1}{2}} \frac{\mathbf{m}^{(g+1)} - \mathbf{m}^{(g)}}{\sigma^{(g)}} \quad (31)$$

where

$\mathbf{p}_\sigma^{(g)} \in \mathbb{R}^n$  is the conjugate evolution path at generation  $g$ .

$c_\sigma < 1$ . Again,  $1/c_\sigma$  is the backward time horizon of the evolution path (compare (20)). For small  $\mu_{\text{eff}}$ , a time horizon between  $\sqrt{n}$  and  $n$  is reasonable.

$\sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}}$  is a normalization constant, see (24).

$\mathbf{C}^{(g)-\frac{1}{2}} \stackrel{\text{def}}{=} \mathbf{B}^{(g)} \mathbf{D}^{(g)-1} \mathbf{B}^{(g)\top}$ , where  $\mathbf{C}^{(g)} = \mathbf{B}^{(g)} (\mathbf{D}^{(g)})^2 \mathbf{B}^{(g)\top}$  is an eigendecomposition of  $\mathbf{C}^{(g)}$ , where  $\mathbf{B}^{(g)}$  is an orthonormal basis of eigenvectors, and the diagonal elements of the diagonal matrix  $\mathbf{D}^{(g)}$  are square roots of the corresponding positive eigenvalues (cf. Sect. 0.1).



# CMA-ES

## Step size[6]

$$\sigma^{(g+1)} = \sigma^{(g)} \exp \left( \frac{c_\sigma}{d_\sigma} \left( \frac{\|\mathbf{p}_\sigma^{(g+1)}\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)$$

Final sampling rule:

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

## CMA-ES extensions

## CMA-ES extensions

Main ideas:

- Parameter modification
- Covariance matrix adaptation modification (e.g. including worst samples)
- **Restarts**
- **Increasing population size**
- **Surrogate models**

# IPOP-CMA-ES[1]

## Restart criteria:

- Stop if the range of the best objective function values of the last  $10 + \lceil 30n/ \rceil$  generations is zero (equalfunvalhist), or the range of these function values and all function values of the recent generation is below  $\text{TolFun} = 10^{-12}$ .
- Stop if the standard deviation of the normal distribution is smaller than  $\text{TolX}$  in all coordinates and  $\vec{p}_c$  (the evolution path from Eq. 2 in [3]) is smaller than  $\text{TolX}$  in all components. We set  $\text{TolX} = 10^{-12}$  (0).
- Stop if adding a 0.1-standard deviation vector in a principal axis direction of  $C^{(g)}$  does not change  $\langle \vec{x}_W^{(g)} \rangle$  (noeffectaxis)<sup>3</sup>
- Stop if adding 0.2-standard deviation in each coordinate does change  $\langle \vec{x}_W^{(g)} \rangle$  (noeffectcoord).
- Stop if the condition number of the covariance matrix exceeds  $10^{14}$  (conditioncov).

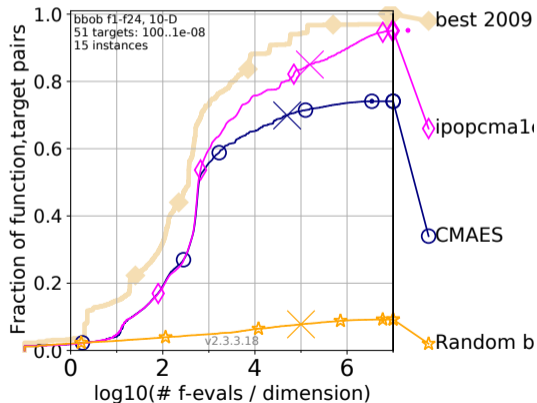
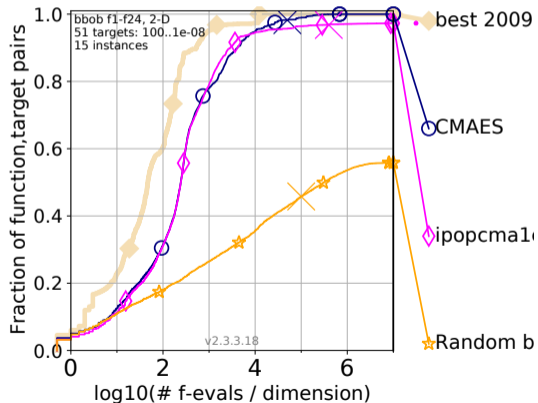
## A Restart CMA Evolution Strategy With Increasing Population Size

**Anne Auger**  
CoLab Computational Laboratory,  
ETH Zürich, Switzerland  
Anne.Auger@inf.ethz.ch

**Nikolaus Hansen**  
CSE Lab,  
ETH Zürich, Switzerland  
Nikolaus.Hansen@inf.ethz.ch

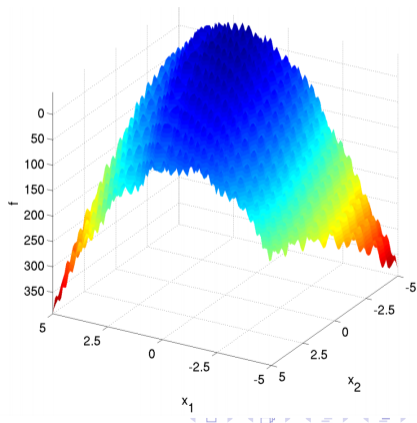
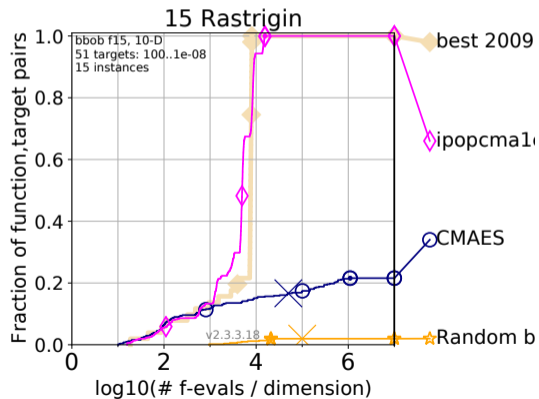
# IPOP-CMA-ES[1]

## Comparison



# IPOP-CMA-ES[1]

## Comparison



## LQ-CMA-ES[7]

Key concept:

- Linear / quadratic surrogate model added to the CMA-ES
- The purpose of the model is to reduce the search cost (true evaluations)
- Weighted linear regression is used to obtain coefficients
- Rank correlation is used as meta-model quality measure
- Model optimum is injected to the population in the next iteration

### A Global Surrogate Assisted CMA-ES

Nikolaus Hansen  
Inria & Ecole polytechnique  
Palaiseau, France  
forename.lastname@inria.fr

# LQ-CMA-ES[7]

---

## Algorithm 1 Determine Population Surrogate Values

---

**Require:** A population  $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_\lambda$ , a model  $\mathcal{M}$  with a data queue of at most  $\max(\lambda, 2df_{\max})$  pairs  $(\mathbf{y}_i, f(\mathbf{y}_i))$ , and a fitness function  $f$

```

1:  $k \leftarrow \lceil 1 + \max(\lambda \times 2\%, 3/0.75 - |\mathcal{M}|) \rceil$  # incrementing evaluations
2: while  $|\mathbf{X}| > 0$  do # while not all elements are added to  $\mathcal{M}$ 
3:   drop the  $k - (\lambda - |\mathbf{X}|)$   $\mathcal{M}$ -best elements from  $\mathbf{X}$  into  $\mathcal{M}$ 
4:   sort the newest  $\min(k, \lambda)$  elements in  $\mathcal{M}$  w.r.t.  $f$  # last = best
5:    $\mathbf{y}_1, \dots, \mathbf{y}_j \leftarrow$  the last  $\max(15, \min(1.2k, 0.75\lambda))$  elements in  $\mathcal{M}$ 
6:   if Kendall- $\tau([\mathcal{M}(\mathbf{y}_i)]_i, [f(\mathbf{y}_i)]_i) \geq 0.85$  then
7:     break while
8:    $k \leftarrow \lceil 1.5k \rceil$ 
9: if  $|\mathbf{X}| > 0$  then
10:  return  $\mathcal{M}(\mathbf{x}_1), \dots, \mathcal{M}(\mathbf{x}_\lambda)$  all offset by
11:     $\min_{\mathbf{x} \in \text{last } k \text{ elements of } \mathcal{M}}(f(\mathbf{x})) - \min_{i=1 \dots \lambda}(\mathcal{M}(\mathbf{x}_i))$ 
12: else
13:  return  $f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda)$ 

```

---

$$z_{\text{lin}} : \mathbf{x} \mapsto [1, x_1, x_2, \dots, x_n]^\top$$

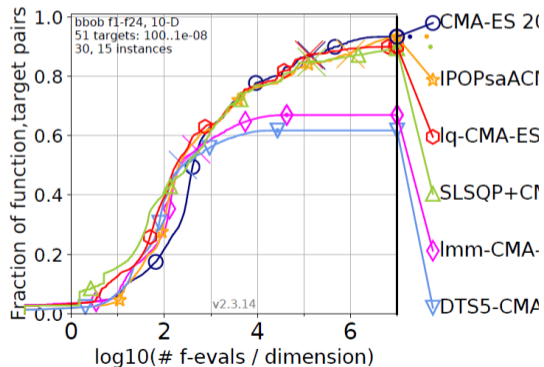
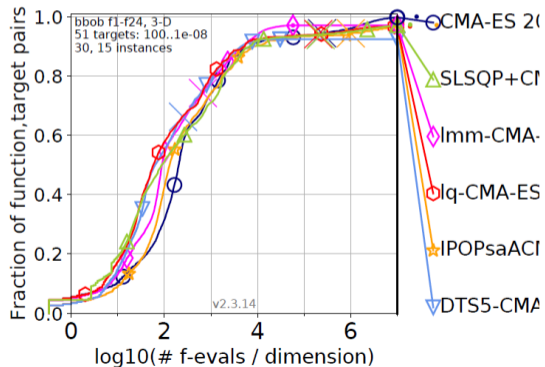
$$z_{\text{quad}} : \mathbf{x} \mapsto [z_{\text{lin}}(\mathbf{x})^\top, x_1^2, \dots, x_n^2]^\top$$

$$z_{\text{full}} : \mathbf{x} \mapsto [z_{\text{quad}}(\mathbf{x})^\top, x_1x_2, x_1x_3, \dots, x_1x_n, \\ x_2x_3, \dots, x_2x_n, x_3x_4, \dots, x_{n-1}x_n]^\top$$



# LQ-CMA-ES[7]

## Comparison



# Experiments

# CMA-ES + Quadratic model + rotation + $R^2$

## Concept

- CMA-ES-based with simple restart criterion and population doubling
- Add FIFO archive evaluated samples (size:  $5 * \lambda$ )
- Extend (conditionally) population with  $\lambda + 1$  offspring
- $R^2$  criterion
- $\lambda + 1$  offspring is designated by the hyper-parabola peak
- Transform archive population using  $B^T$  matrix from SVD decomposition of covariance matrix  $C$

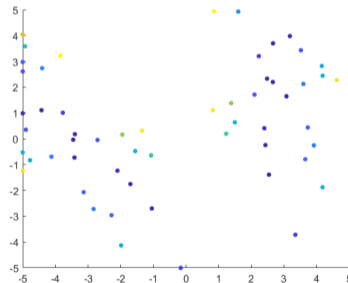
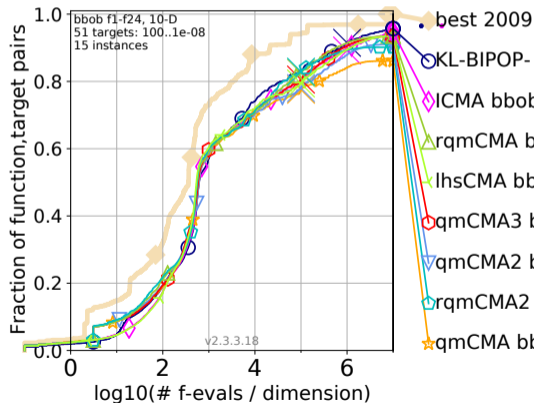
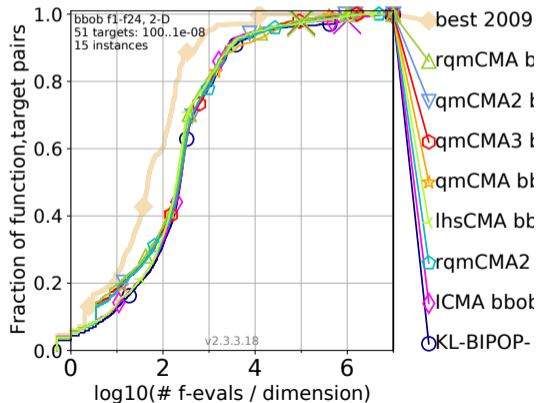


Figure 13: Step Ellipsoidal example

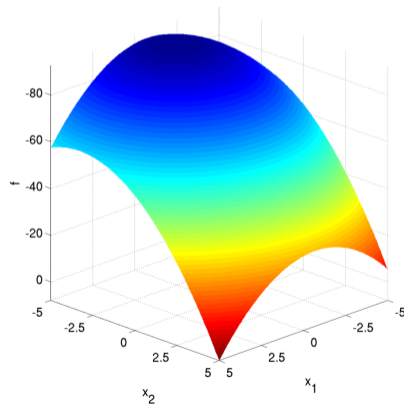
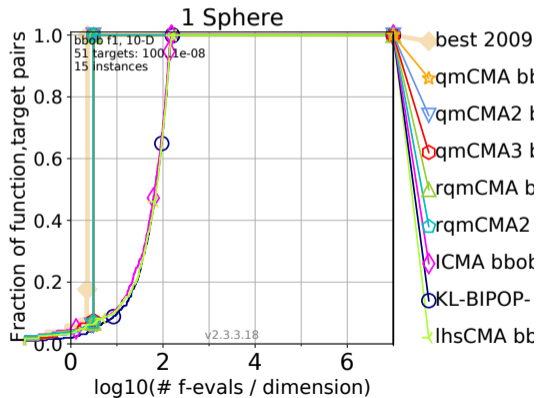
# CMA-ES + Quadratic model + rotation + $R^2$

## Comparison



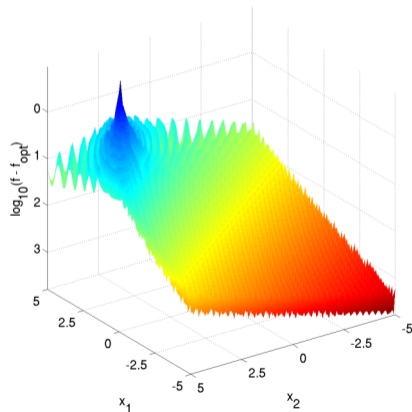
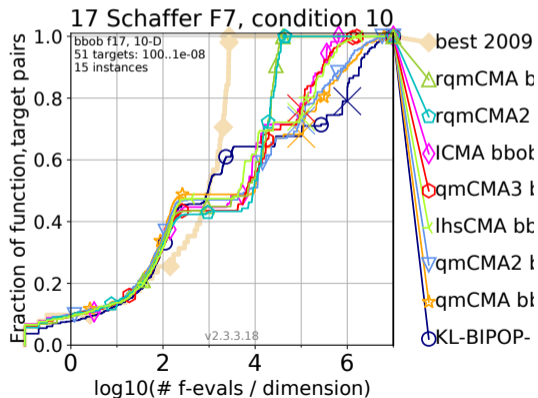
# CMA-ES + Quadratic model + rotation + $R^2$

## Comparison



# CMA-ES + Quadratic model + rotation + $R^2$

## Comparison



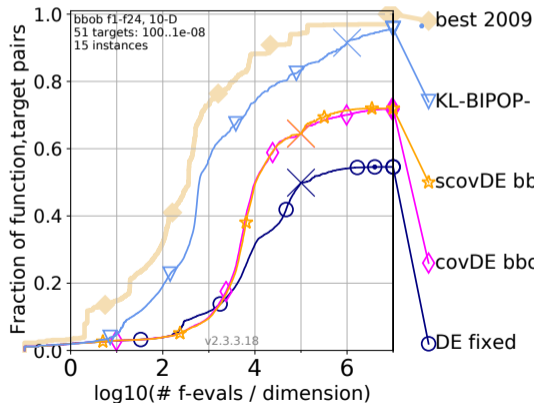
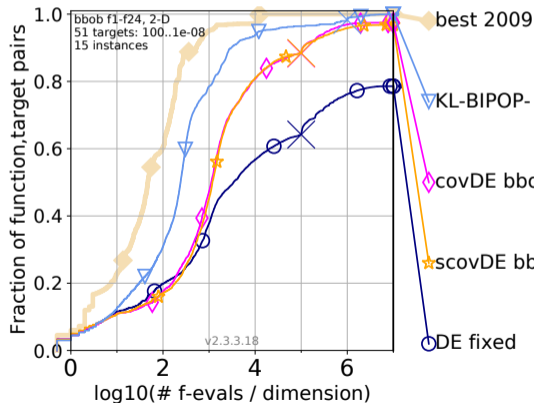
## DE + rotation

### Concept:

- Differential Evolution based with simple restart criterion
- Covariance matrix  $C$  calculated after each iteration using  $\lambda$  offspring (current population)
- Covariance matrix  $C$  arithmetic smoothing is possible ( $0.2 * C^g$ )
- Covariance matrix  $C$  SVD to get transformation (rotation) component
- Transformation is done before mutation
- Reverse-transformation is done before selection

# DE + rotation

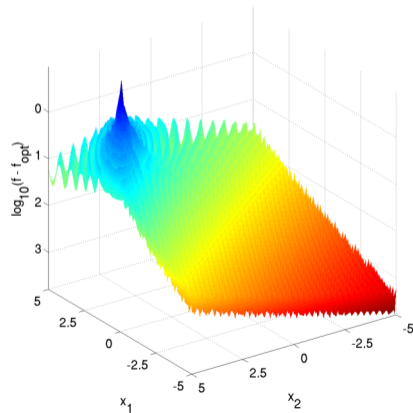
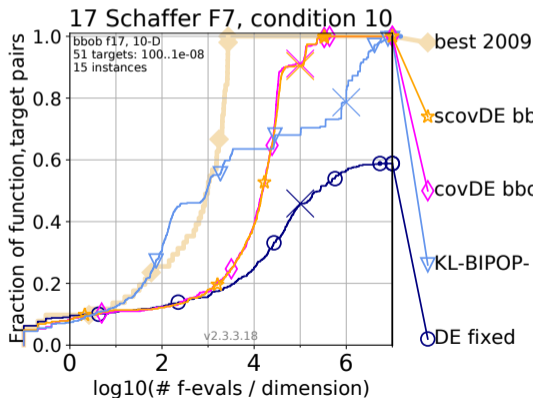
## Comparison





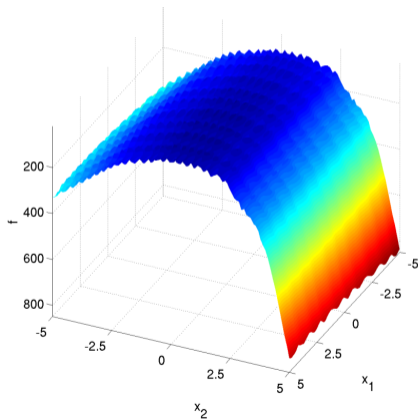
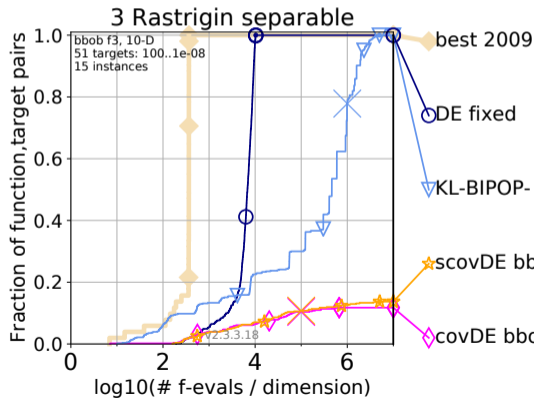
# DE + rotation

## Comparison



# DE + rotation










## Comparison



## Summary

- CMA-ES is a powerful optimization algorithm
- CMA-ES behaves well on strong global structure functions
- CMA-ES can designate optimum of rotated hyper-parabola by its normal sampling procedure
- Rotation concept can be applied into Differential Evolution
- State-of-the-art version of DE such as L-SHADE should enhanced with rotation component

## Bibliography

-  Auger, A., Hansen, N.: A restart cma evolution strategy with increasing population size. In: 2005 IEEE congress on evolutionary computation. vol. 2, pp. 1769–1776. IEEE (2005)
-  Auger, A., Hansen, N.: Tutorial cma-es: evolution strategies and covariance matrix adaptation. In: Proceedings of the 14th annual conference companion on Genetic and evolutionary computation. pp. 827–848 (2012)
-  Commons, W.: Concept of directional optimization in cma-es algorithm (2020), [https://commons.wikimedia.org/wiki/File:Concept\\_of\\_directional\\_optimization\\_in\\_CMA-ES\\_algorithm.png](https://commons.wikimedia.org/wiki/File:Concept_of_directional_optimization_in_CMA-ES_algorithm.png)
-  Commons, W.: Illustration of a multivariate gaussian distribution and its marginals. (2020), [https://en.wikipedia.org/wiki/Multivariate\\_normal\\_distribution#/media/File:MultivariateNormal.png](https://en.wikipedia.org/wiki/Multivariate_normal_distribution#/media/File:MultivariateNormal.png)
-  Hansen, N.: The cma evolution strategy: a comparing review. In: Towards a new evolutionary computation, pp. 75–102. Springer (2006)
-  Hansen, N.: The cma evolution strategy: A tutorial. arXiv preprint arXiv:1604.00772 (2016)
-  Hansen, N.: A global surrogate assisted cma-es. In: Proceedings of the Genetic and Evolutionary Computation Conference. pp. 664–672 (2019)
-  Hansen, N., Brockhoff, D., Mersmann, O., Tusar, T., Tusar, D., ElHara, O.A., Sampaio, P.R., Atamna, A., Varelas, K., Batu, U., Nguyen, D.M., Matzner, F., Auger, A.: COmparing Continuous Optimizers: numbbo/COCO on Github (2019). <https://doi.org/10.5281/zenodo.2594848>, <https://doi.org/10.5281/zenodo.2594848>
-  Hansen, N., Ostermeier, A.: Completely derandomized self-adaptation in evolution strategies. Evolutionary computation 9(2), 159–195 (2001)