

# CMA-ES : Une méthode d'optimisation sans dérivées pour des problèmes difficiles

Anne Auger

RandOpt, Inria et CMAP, Ecole Polytechnique

Rencontres de la Chaire Modélisation Mathématique et Biodiversité

6 mai 2025



# Derivative-free / Gradient-free optimization

**Minimize**  $f: x = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto f(x)$

approach  $x^\star$  such that  $f(x^\star) \leq f(x)$  for all  $x$

# Derivative-free / Gradient-free optimization

**Minimize**  $f: x = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto f(x)$

approach  $x^\star$  such that  $f(x^\star) \leq f(x)$  for all  $x$

**Assume:** derivatives or gradient of  $f$  **not available**

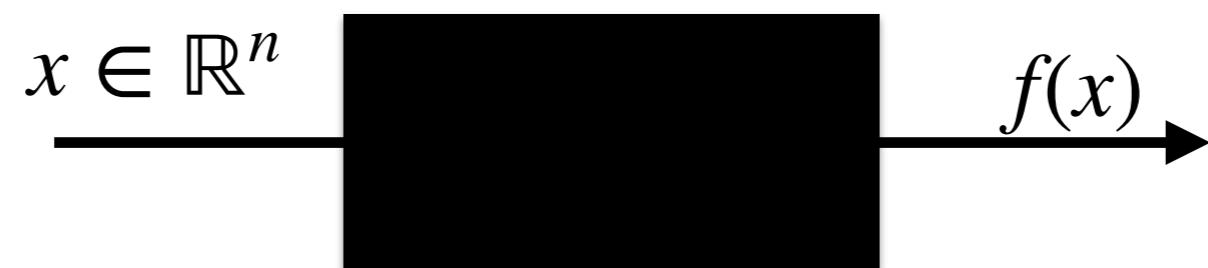
# Derivative-free / Gradient-free optimization

**Minimize**  $f: x = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto f(x)$

approach  $x^\star$  such that  $f(x^\star) \leq f(x)$  for all  $x$

**Assume:** derivatives or gradient of  $f$  **not available**

**Zero-order black-box optimization**



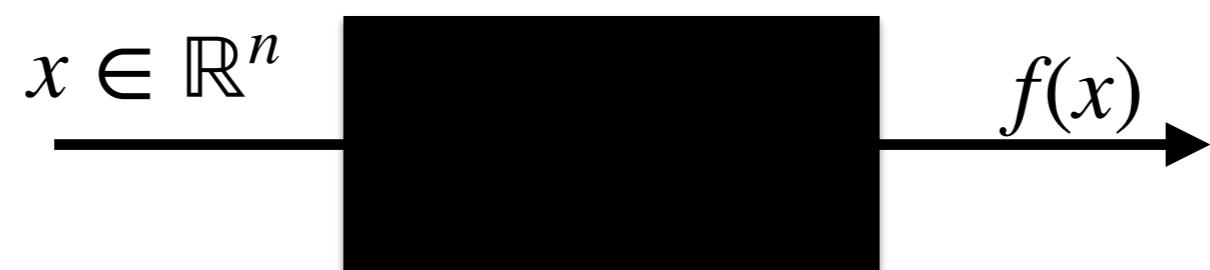
# Derivative-free / Gradient-free optimization

**Minimize**  $f: x = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto f(x)$

approach  $x^\star$  such that  $f(x^\star) \leq f(x)$  for all  $x$

**Assume:** derivatives or gradient of  $f$  **not available**

**Zero-order black-box optimization**



**Many applications**

(see also presentation Renaud)

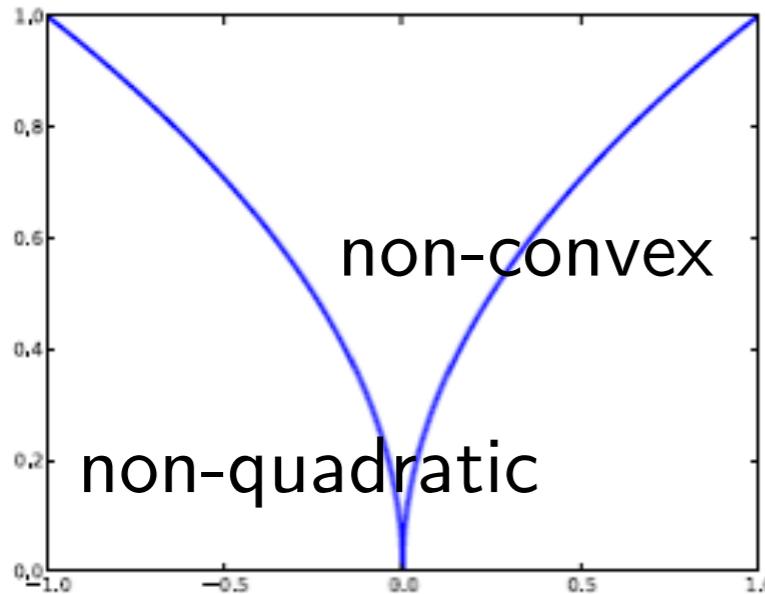


# Why stochastic numerical black-box optimization?

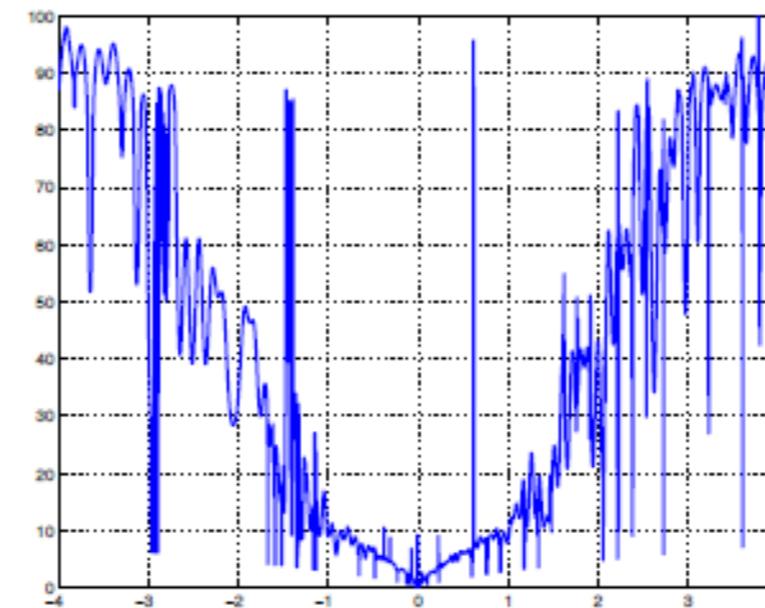
Difficulties inside the black-box

?

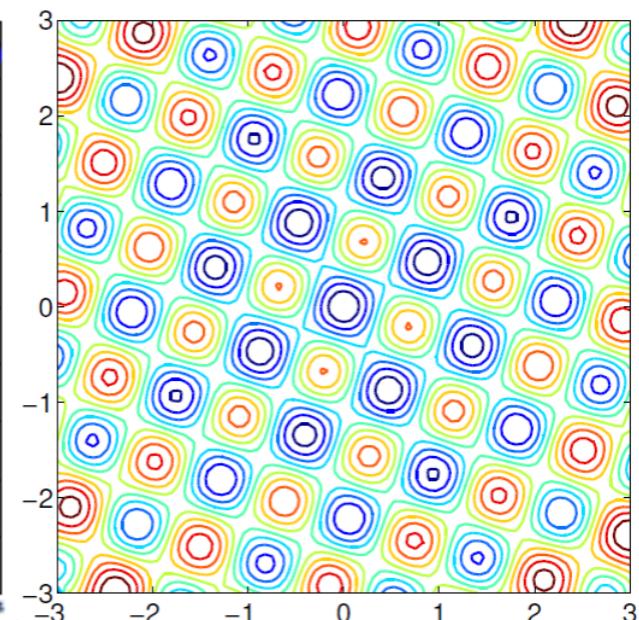
non-linear



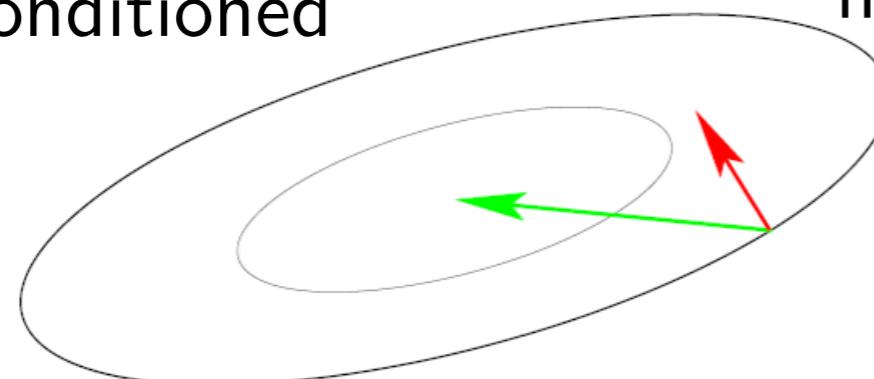
non-differentiable



non-separable



ill-conditioned



noisy

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

# From Evolution Strategies (ES) to CMA-ES

Evolution Strategies (70s+)  
randomized bio-inspired population-based algorithms



Rechenberg Schwefel

# From Evolution Strategies (ES) to CMA-ES

Evolution Strategies (70s+)  
randomized bio-inspired population-based algorithms



Rechenberg Schwefel

CMA-ES (1996/2001 - today)

More than **70 million downloads** of its two main Python implementations

# From Evolution Strategies (ES) to CMA-ES

Evolution Strategies (70s+)  
randomized bio-inspired population-based algorithms



Rechenberg Schwenefel

CMA-ES (1996/2001 - today)  
not enough bio-inspired?

The first fathers [in Rechenberg's lab in Berlin ]



Hansen

Ostermeier

More than **70 million downloads** of its two main Python implementations

# From Evolution Strategies (ES) to CMA-ES

Evolution Strategies (70s+)  
randomized bio-inspired population-based algorithms



Rechenberg Schwenefel

CMA-ES (1996/2001 - today)  
not enough bio-inspired?

The first fathers [in Rechenberg's lab in Berlin ]

RandOpt: Auger/Hansen/Brockhoff+Chotard



Akimoto

Glasmachers

Hansen

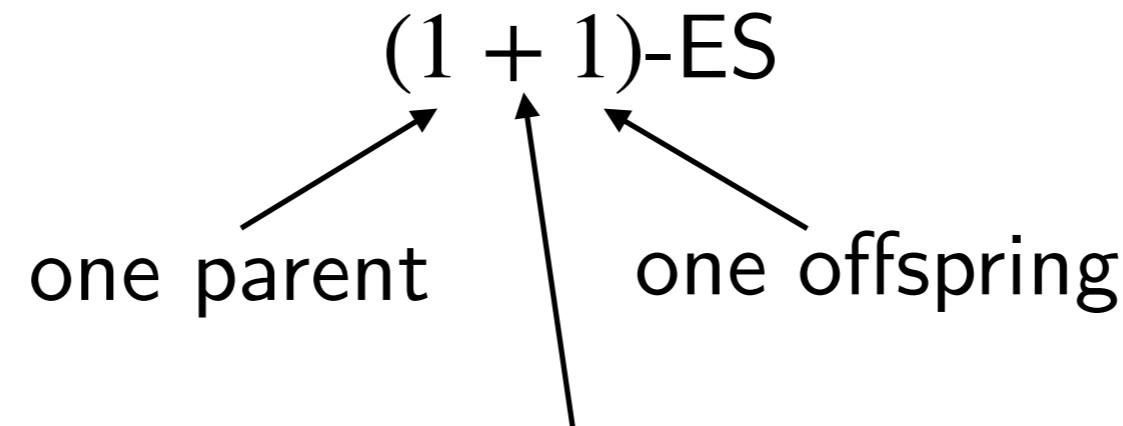
Ostermeier

Arnold

More than **70 million downloads** of its two main Python implementations

From a simple ES to CMA-ES

# A simple Evolution Strategy: (1+1)-ES



Elitist selection: keep best among  
parent and offspring

# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

# A simple Evolution Strategy: (1+1)-ES

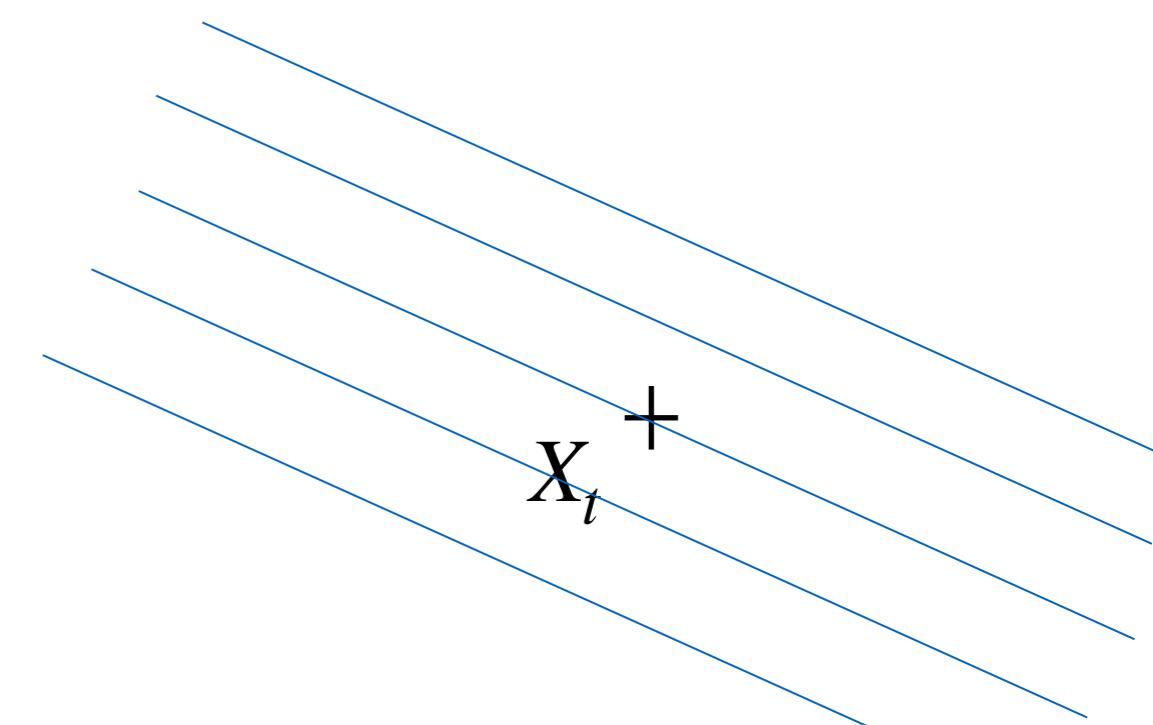
$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

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

$\text{minimize } f: \mathbb{R}^2 \rightarrow \mathbb{R}$



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

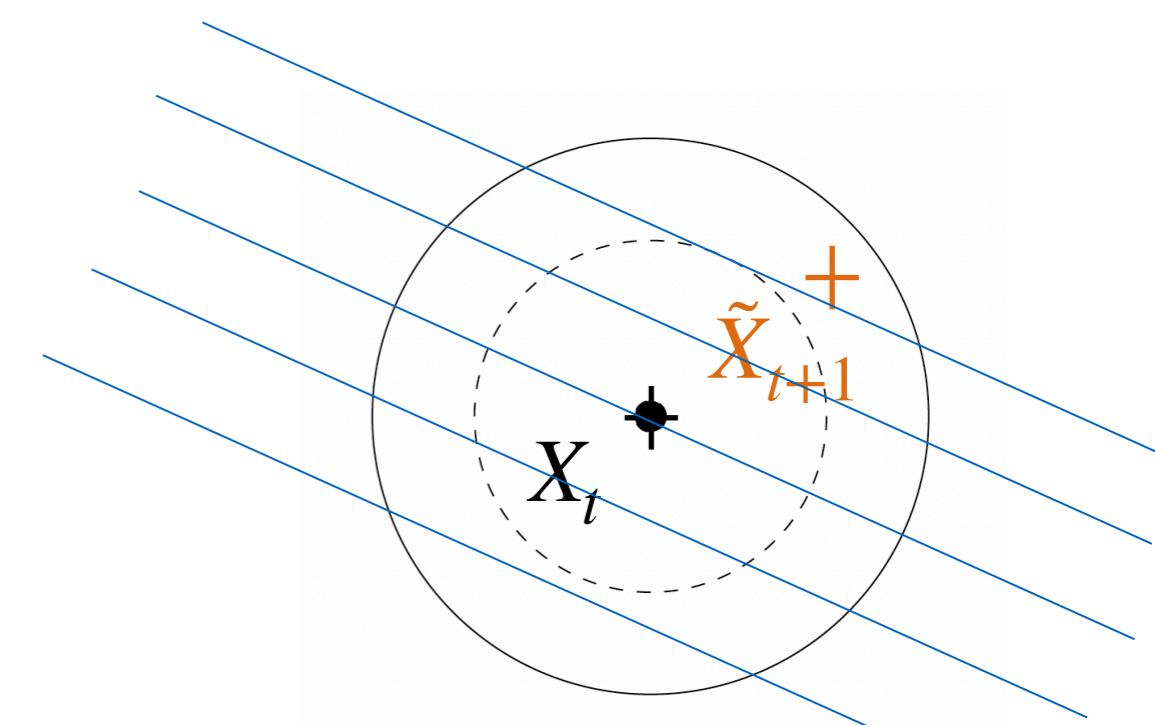
$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

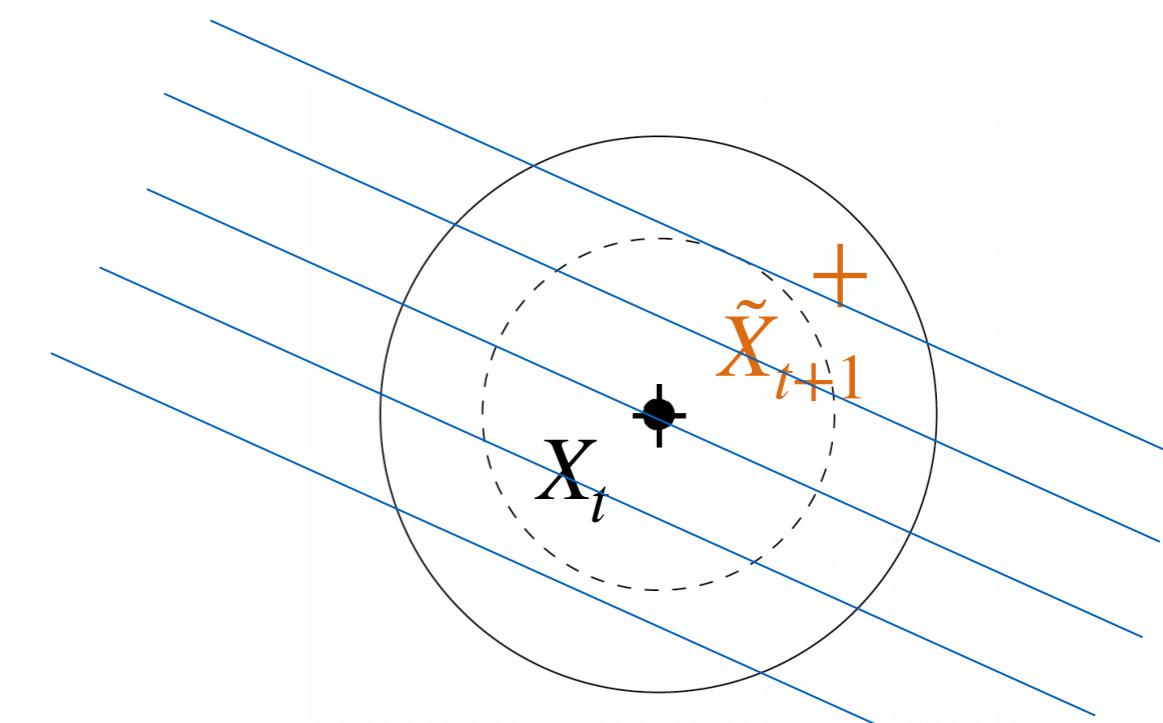
$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

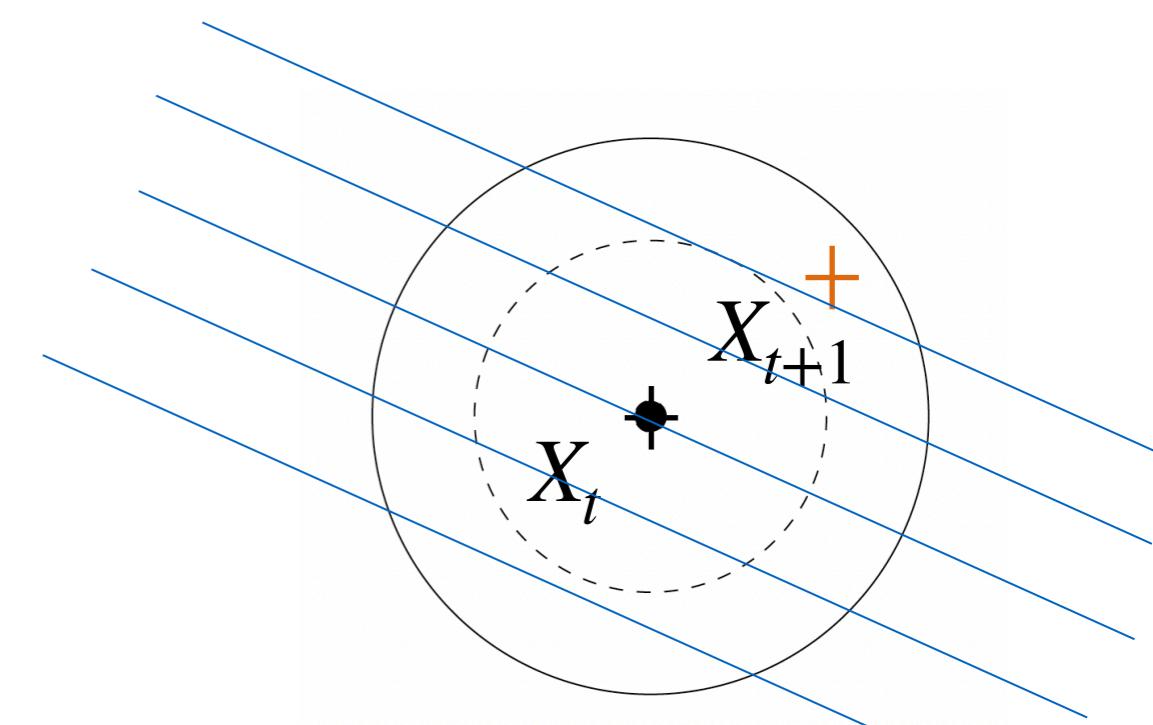
$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

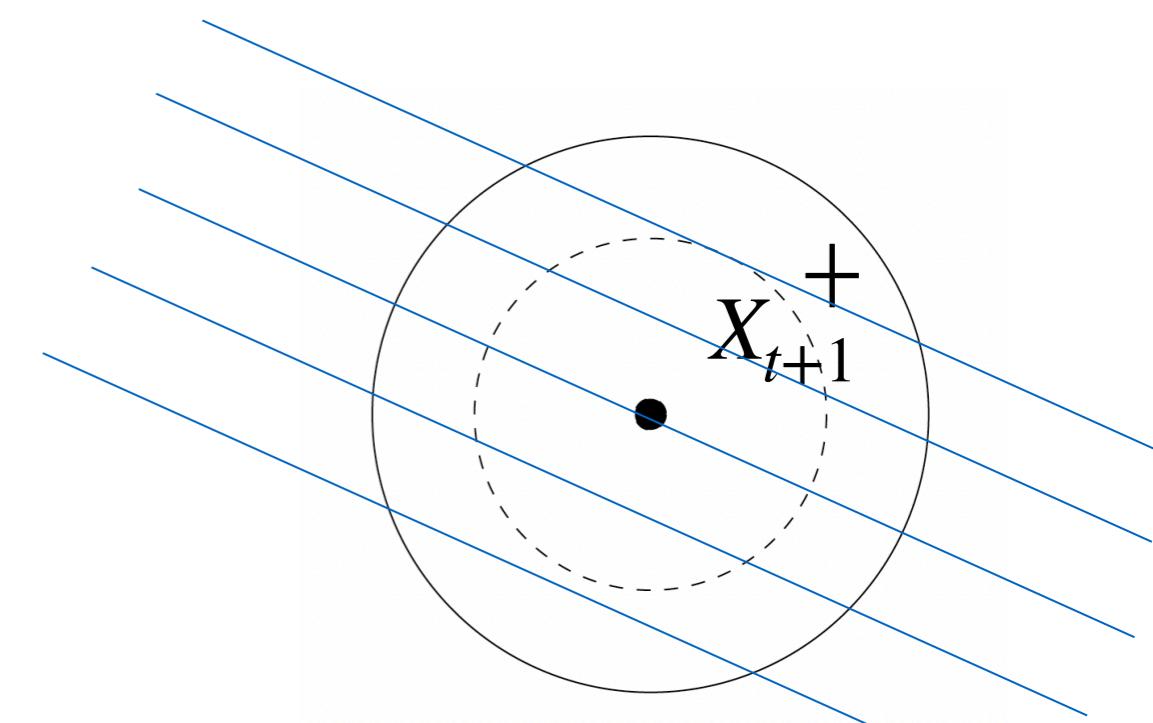
$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

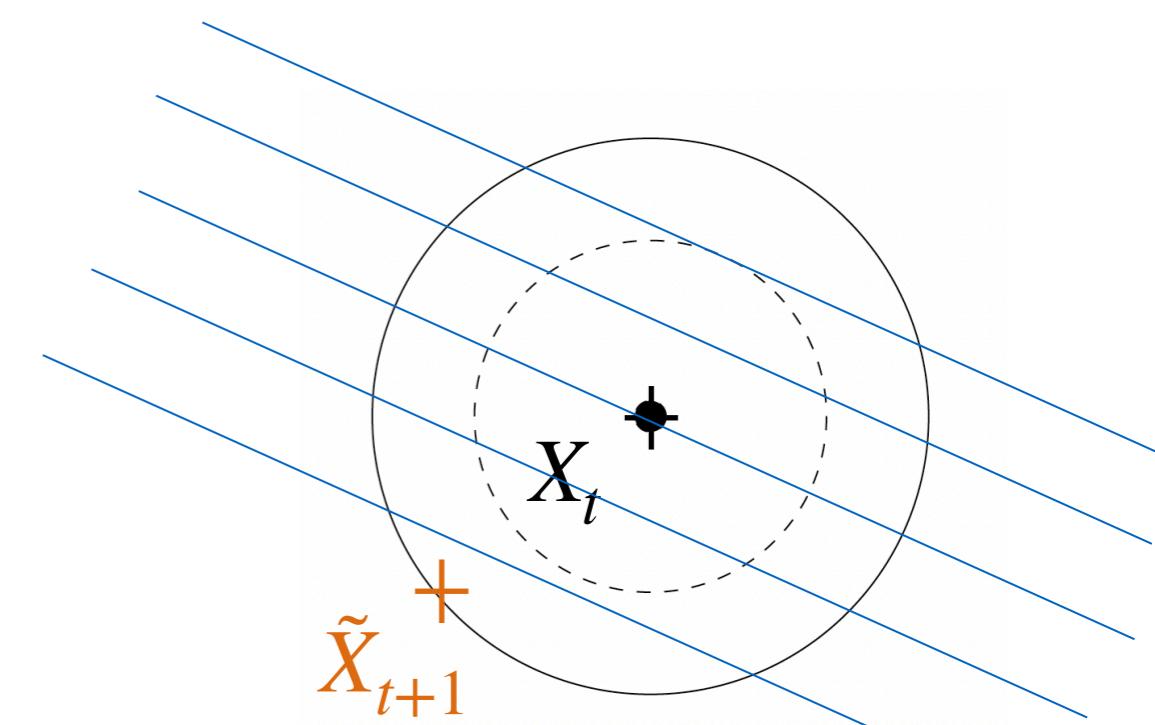
$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

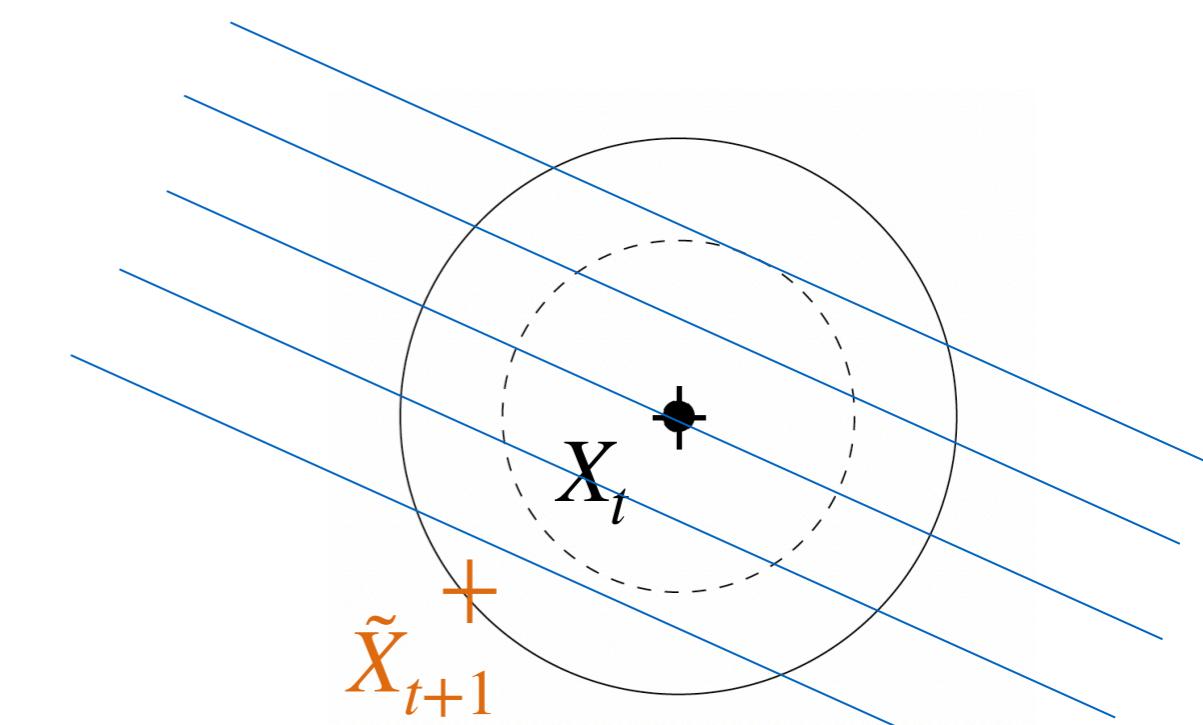
IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring

ELSE

$X_{t+1} = X_t$  #keep current parent



# A simple Evolution Strategy: (1+1)-ES

$t \in \mathbb{N}$  iteration index

$X_t \in \mathbb{R}^n$  : parent at iteration t

$\sigma \in \mathbb{R}_>$  : step-size

$\tilde{X}_{t+1}$  : offspring created via **mutation** of parent

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

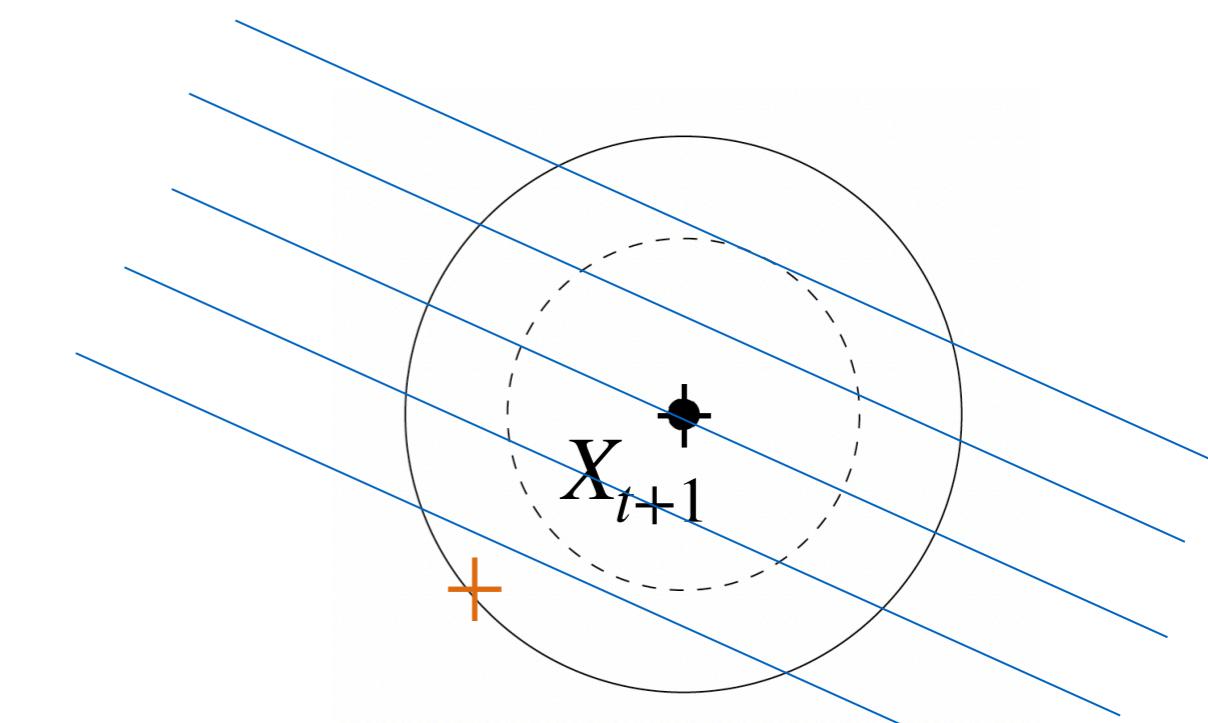
IF offspring better than parent [  $f(\tilde{X}_{t+1}) \leq f(X_t)$  ]

minimize  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring

ELSE

$X_{t+1} = X_t$  #keep current parent



(1+1)-ES algorithm - minimize  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

WHILE not happy

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF  $f(\tilde{X}_{t+1}) \leq f(X_t)$  THEN

$X_{t+1} = \tilde{X}_{t+1}$  #keep offspring

ELSE

$X_{t+1} = X_t$  #keep current parent

ENDIF

$$t \leftarrow t + 1$$

(1+1)-ES algorithm - minimize  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

WHILE not happy

$$\tilde{X}_{t+1} = X_t + \sigma \mathcal{N}(0, I_d)$$

IF  $f(\tilde{X}) < f(X_t)$  THEN

Issue #1:

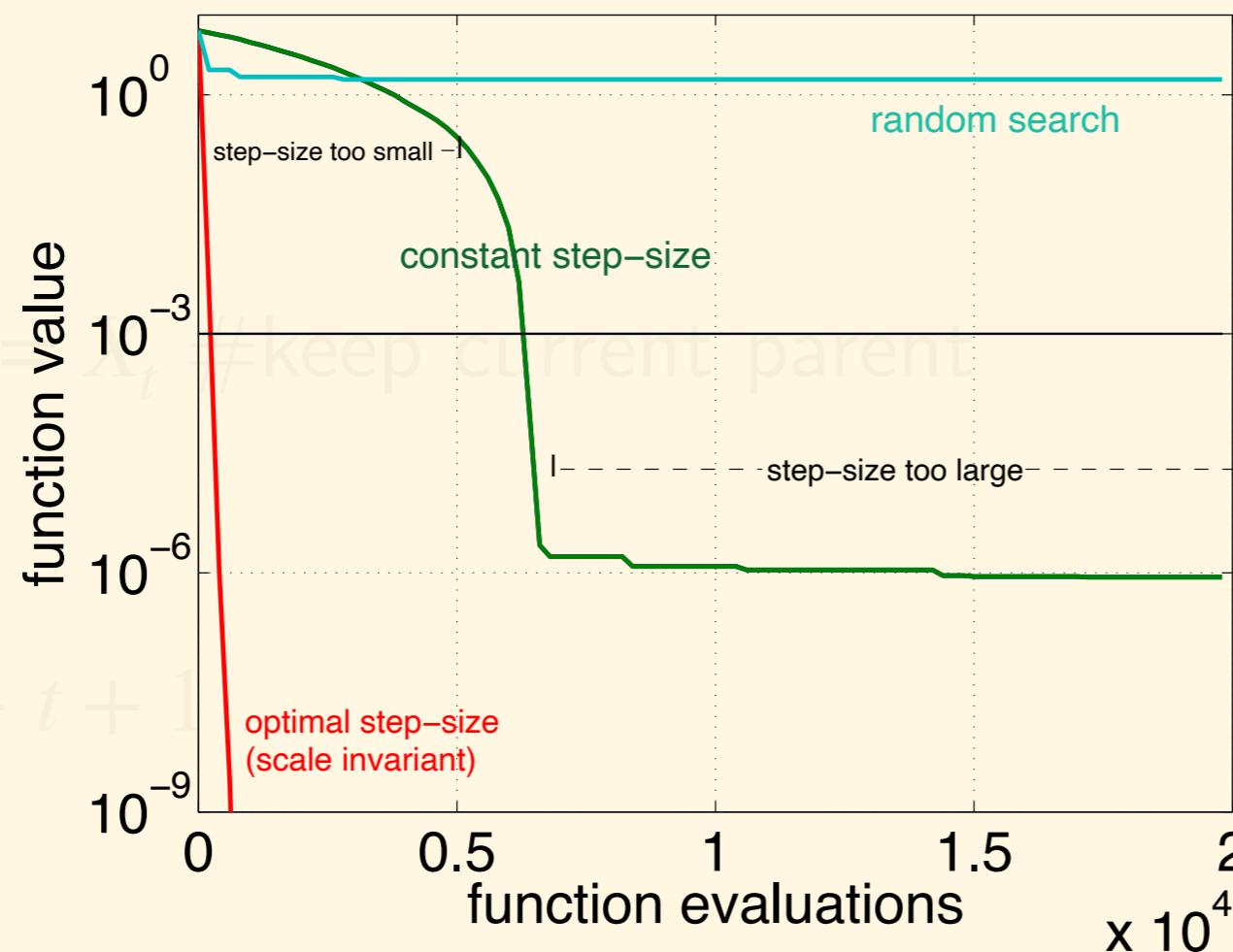
- cannot converge “fast” if  $\sigma$  fixed

ELSE

$X_{t+1} = X_t$

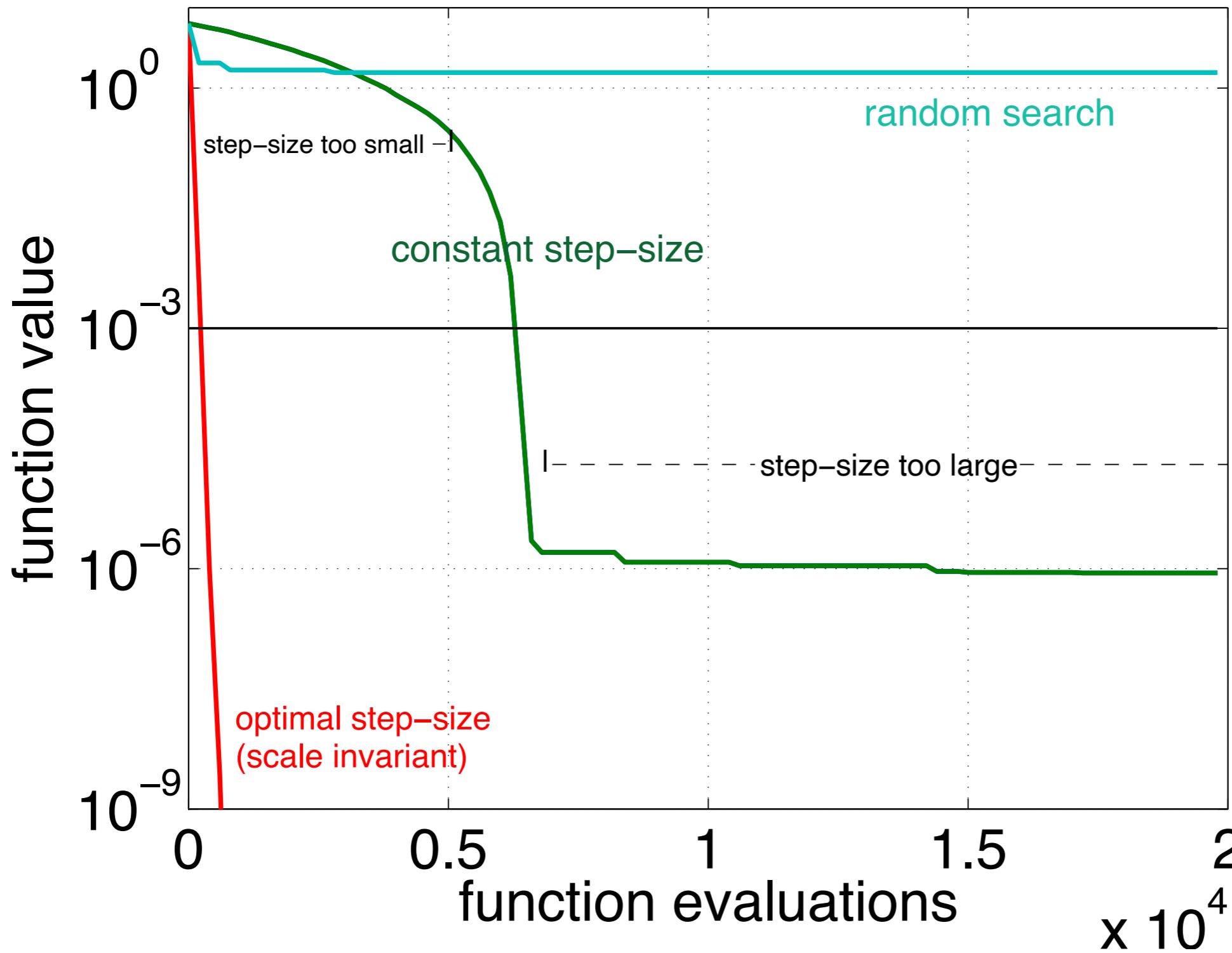
ENDIF

$t \leftarrow t + 1$



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

n=10



Need to adapt the step-size  $\rightarrow$  step-size adaptive ES

# (1+1)-ES algorithm with 1/5-th success rule [Rechenberg]

WHILE not happy

$$\tilde{X}_{t+1} = X_t + \sigma_t \mathcal{N}(0, I_d)$$

IF  $f(\tilde{X}_{t+1}) \leq f(X_t)$  THEN

$$X_{t+1} = \tilde{X}_{t+1} \text{ #keep offspring}$$

$$\sigma_{t+1} = 1.5 \times \sigma_t \text{ #increase step-size}$$

ELSE

$$X_{t+1} = X_t \text{ #keep current parent}$$

$$\sigma_{t+1} = 1.5^{-1/4} \times \sigma_t \text{ #decrease step-size}$$

ENDIF

$$t \leftarrow t + 1$$

# (1+1)-ES algorithm with 1/5-th success rule [Rechenberg]

WHILE not happy

$$\tilde{X}_{t+1} = X_t + \sigma_t \mathcal{N}(0, I_d)$$

IF  $f(\tilde{X}_{t+1}) \leq f(X_t)$  THEN

$$X_{t+1} = \tilde{X}_{t+1}$$

$$\sigma_{t+1} = 1.5 \times \sigma_t$$

ELSE

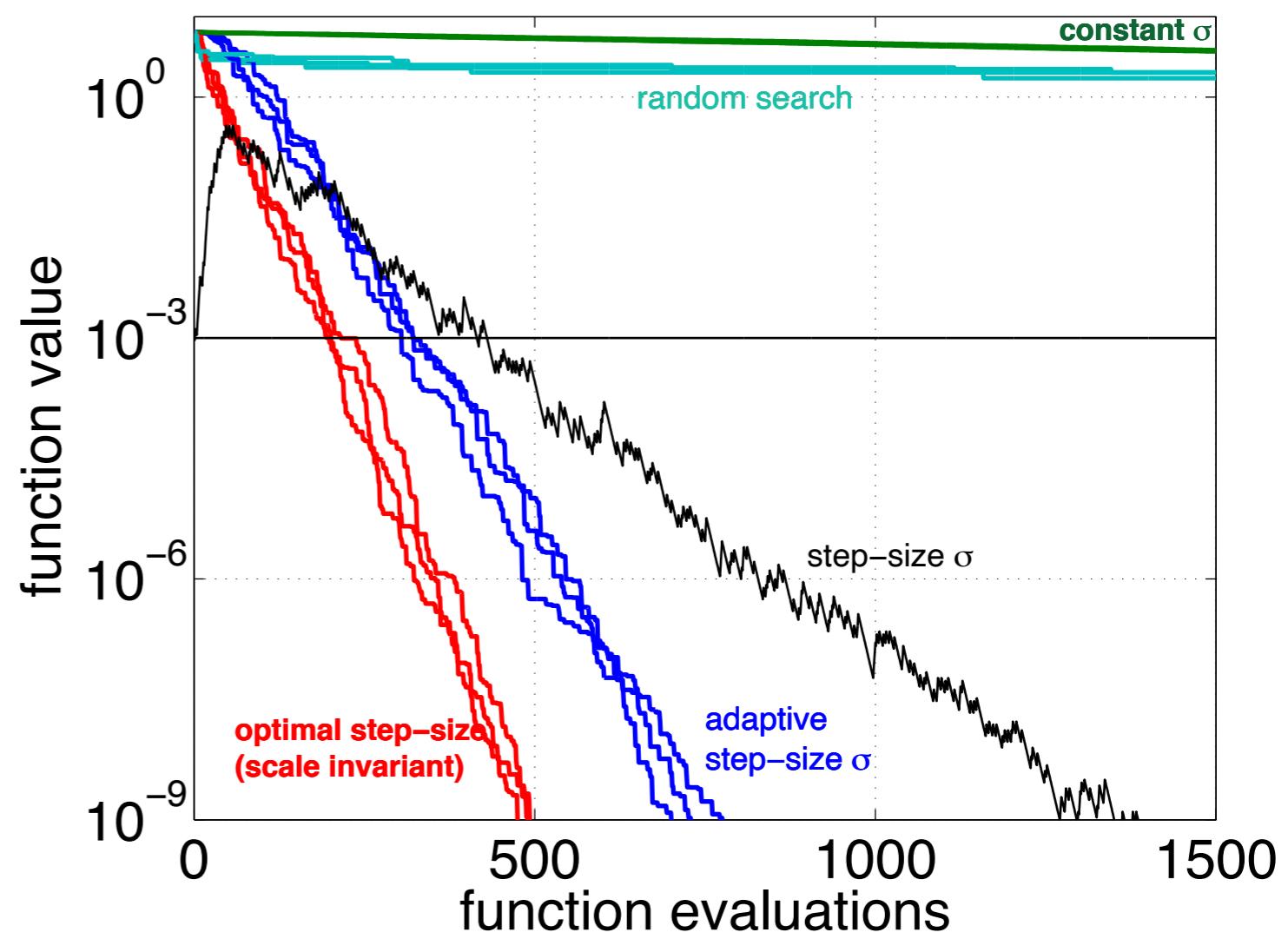
$$X_{t+1} = X_t$$

$$\sigma_{t+1} = 1.5^{-1/4} \times \sigma_t$$

ENDIF

$$t \leftarrow t + 1$$

Achieves **geometric** (“linear”) convergence



$$f(x) = \sum_{i=1}^n x_i^2, n=10$$

# Invented (independently) in other communities

## Schumer, Steiglitz, 1968 "Adaptive Step Size Random Search"

270

IEEE TRANSACTIONS ON AUTOMATIC CONTROL, VOL. AC-13, NO. 3, JUNE 1968

### Adaptive Step Size Random Search

MICHAEL A. SCHUMER, MEMBER, IEEE, AND KENNETH STEIGLITZ, MEMBER, IEEE

*Abstract*—Fixed step size random search for minimization of functions of several parameters is described and compared with the fixed step size gradient method for a particular surface. A theoretical technique, using the optimum step size at each step, is analyzed. A practical adaptive step size random search algorithm is then proposed, and experimental experience is reported that shows the superiority of random search over other methods for sufficiently high dimension.

#### INTRODUCTION

Rastrigin has compared a fixed step size random search (FSSRS) method with a fixed step size gradient method and concluded that under certain circumstances FSSRS is superior. It is clear, however, that if the step size of the random search method were optimum at each step, even better performance would result. In this paper, a hypothetical random search method that uses the optimum step size at each point will be analyzed for a hyperspherical surface. An adaptive step size

## Devroye 1972 "The compound random search algorithm"

THE COMPOUND RANDOM SEARCH ALGORITHM

Luc P. DEVROYE  
Dept. of Electrical Engineering  
Catholic University of Louvain, Belgium

**ABSTRACT**

The optimization of a function of many variables is investigated. A new algorithm is proposed : the compound random search algorithm (CRSA). This algorithm combines the features of random search and non-random direct search. Each part of the algorithm is discussed in detail. The CRSA is compared with several other direct search methods in many different problems. The results are promising. Some modifications of the basic search scheme make the CRSA particularly useful and efficient.

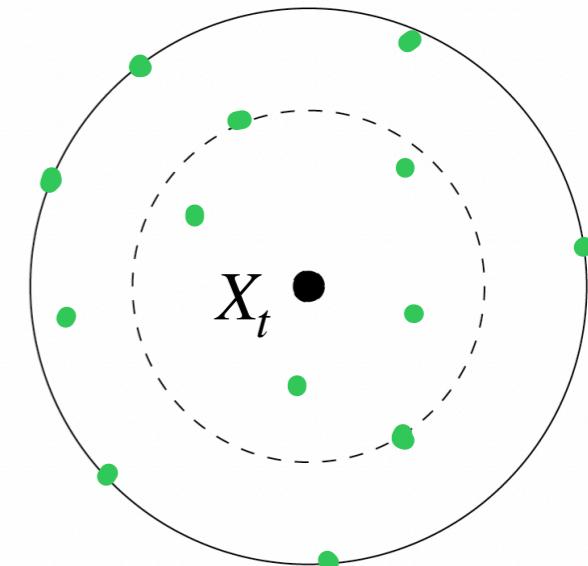
accelerate and control the search. Moreover, new information must be gathered during the search and transformed into useful data. There is a remarkable resemblance between this search and a learning process or a game against nature.

A new algorithm is developed which performs much better, even under the worst circumstances, than the powerful methods of Gucker [13] or Rosenbrock [11]. Moreover, this method is one of the fastest random optimization methods when an optimum must be localized quite accurately. In our discussion, the rate of convergence is determined in

More issues with (1+1)-ES with 1/5th success rule:

candidate solutions are isotropically distributed

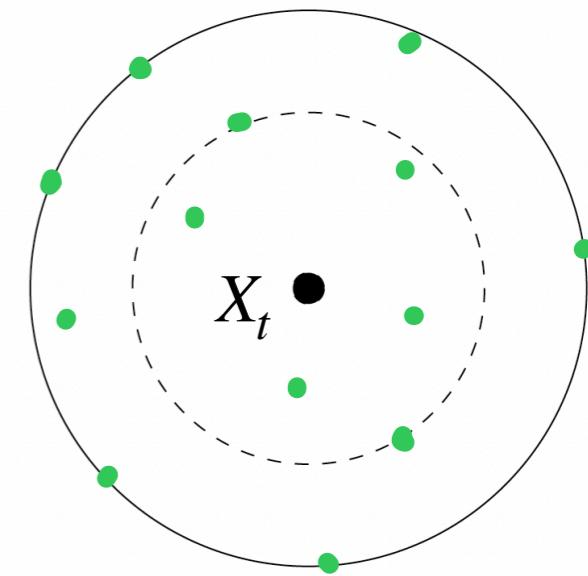
$$\tilde{X}_{t+1} = X_t + \sigma_t \mathcal{N}(0, I_d)$$



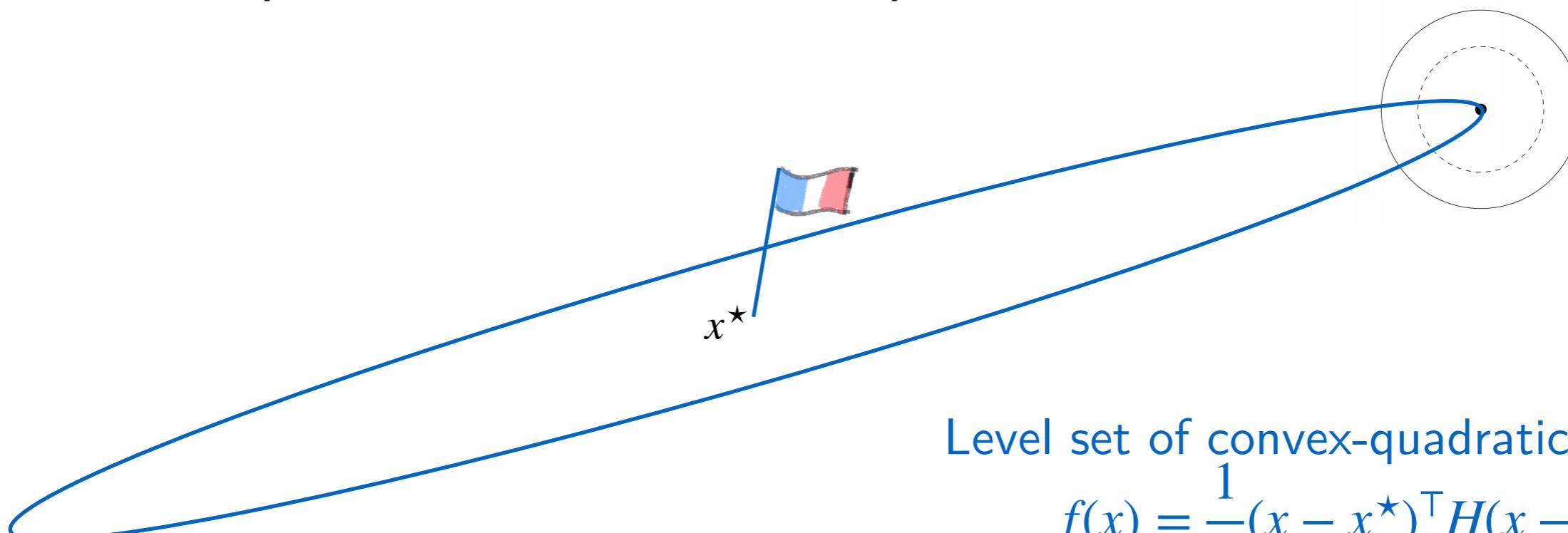
More issues with (1+1)-ES with 1/5th success rule:

candidate solutions are isotropically distributed

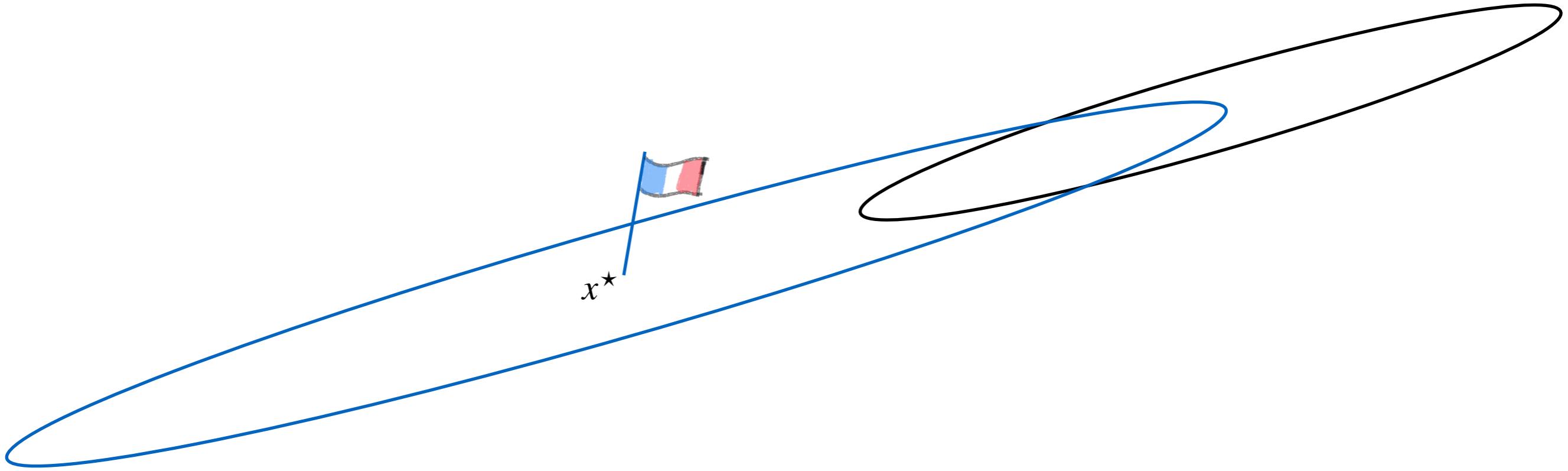
$$\tilde{X}_{t+1} = X_t + \sigma_t \mathcal{N}(0, I_d)$$

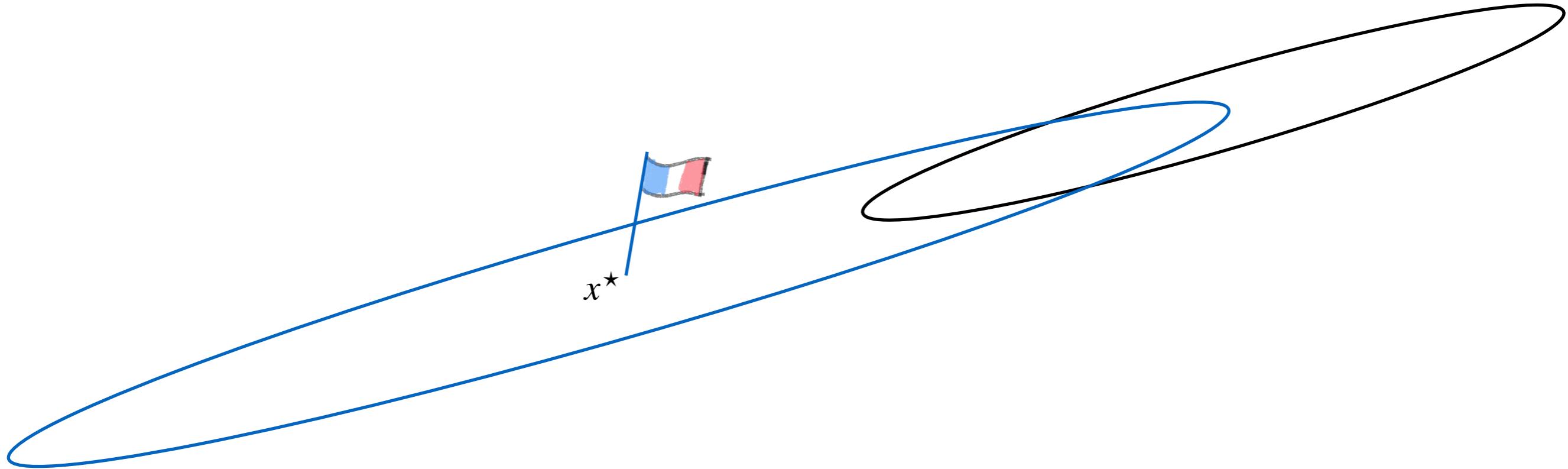


not adapted on ill-conditioned problems:



Level set of convex-quadratic function:  
$$f(x) = \frac{1}{2}(x - x^*)^\top H(x - x^*)$$





Ideally on a ill-conditionned quadratic problem:  $\frac{1}{2}(x - x^*)^\top H(x - x^*)$   
with  $\text{cond}(H) \gg 1$  (order of  $10^6$ ) sample with:

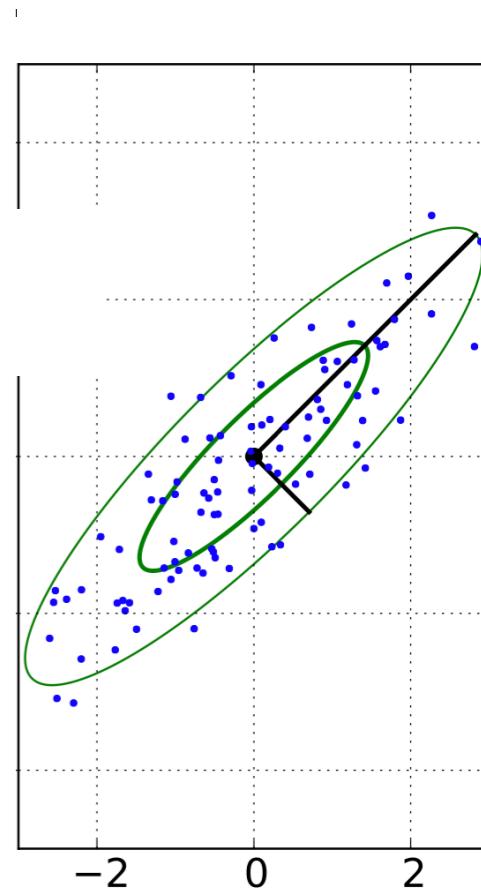
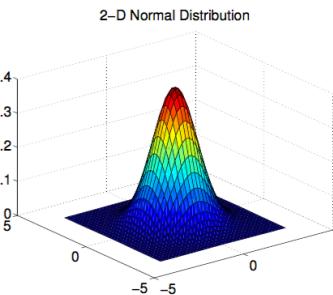
$$X_t + \sigma_t \mathcal{N}(0, C_t) \text{ with } C_t \propto H^{-1}$$

# From (1+1)-ES to $(\mu/\mu_w, \lambda)$ -CMA-ES

- ① adapt the covariance matrix of sampling distribution:  
step-size + covariance matrix adaptive ES = CMA-ES
- ② use a population:  $(\mu/\mu_w, \lambda)$ -CMA-ES
- ③ use mutation and (meta)-crossover

# Adaptive Stochastic Optimization Algorithm

Given e.g.  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$



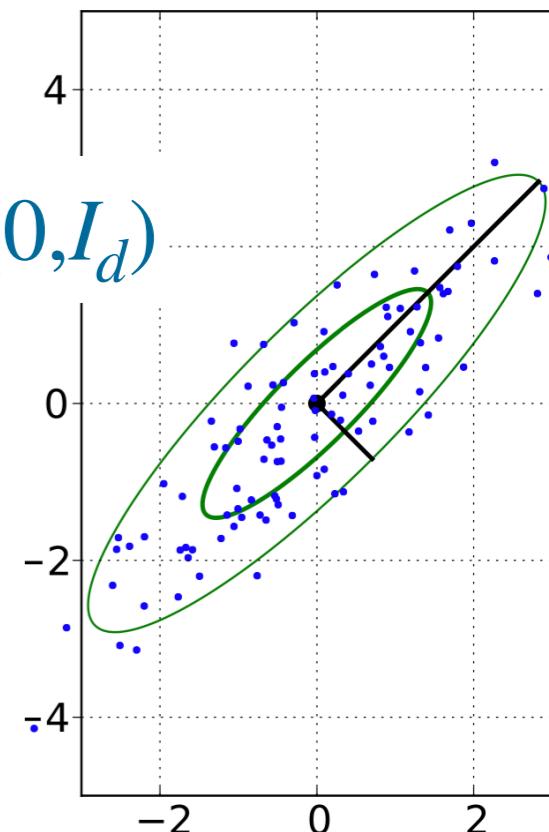
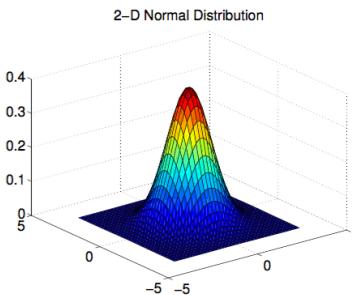
# Adaptive Stochastic Optimization Algorithm

Given e.g.  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$

- ① Sample candidate solutions  $X_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ , i.e.

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i, \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d., } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$



# Adaptive Stochastic Optimization Algorithm

Given e.g.  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$

- ① Sample candidate solutions  $X_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ , i.e.

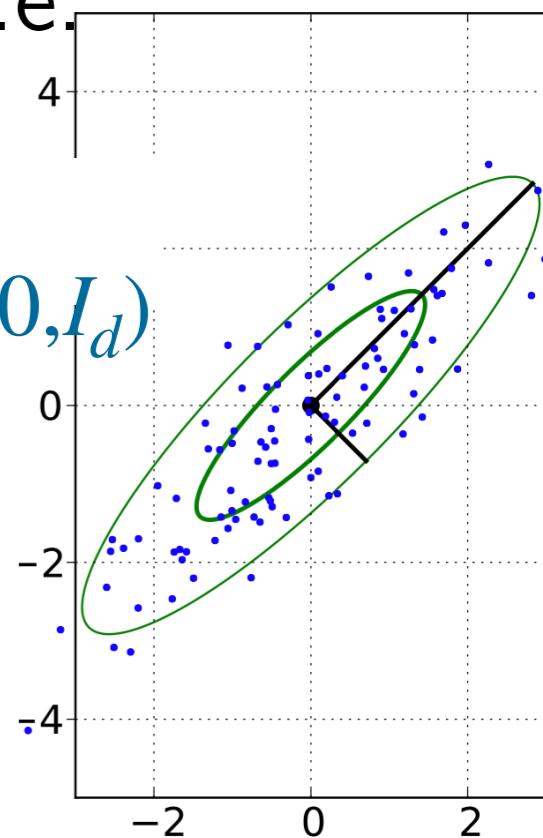
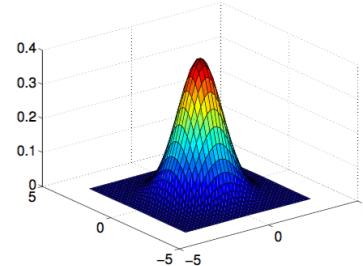
$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i, \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d., } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

- ② Evaluate and rank candidate solutions

$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

2-D Normal Distribution



# Adaptive Stochastic Optimization Algorithm

Given e.g.  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$

- ① Sample candidate solutions  $X_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ , i.e.

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i, \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d., } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

- ② Evaluate and rank candidate solutions

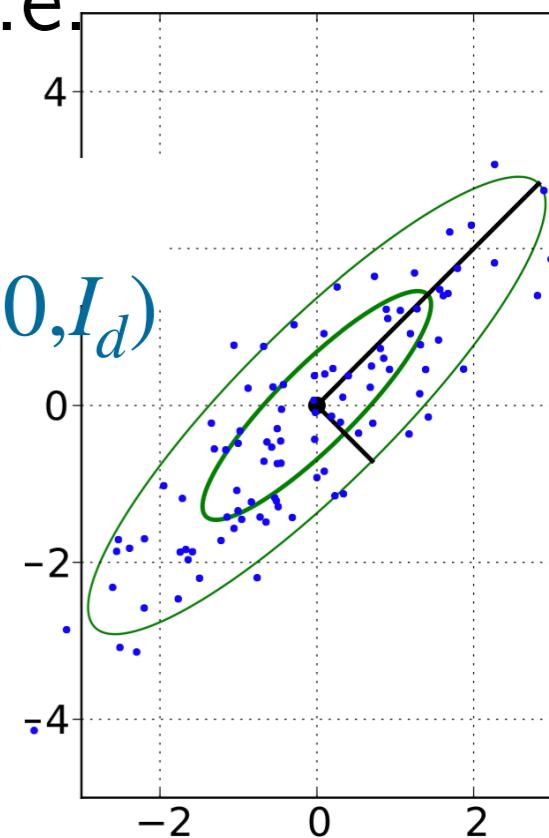
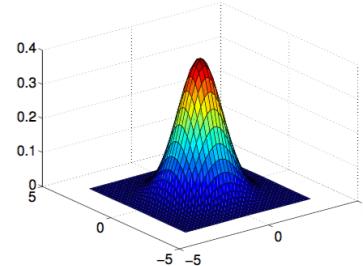
$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

- ③ Update  $\theta_t$ :

$$\theta_{t+1} = G\left(\theta_t, [U_{t+1}^{s_{t+1}(1)}, \dots, U_{t+1}^{s_{t+1}(\lambda)}]\right)$$

*should drive  $m_t$  towards the optimum*

2-D Normal Distribution



## Comparison-based $\Rightarrow$ Invariance strict. increasing functions

Let  $g : \text{Im}f \rightarrow \text{Im}(g)$  be strictly increasing

$$\begin{aligned} f\left(X_{t+1}^{s(1)}\right) &\leq \dots \leq f\left(X_{t+1}^{s(\lambda)}\right) \\ \Leftrightarrow \\ g \circ f\left(X_{t+1}^{s(1)}\right) &\leq \dots \leq g \circ f\left(X_{t+1}^{s(\lambda)}\right) \end{aligned}$$

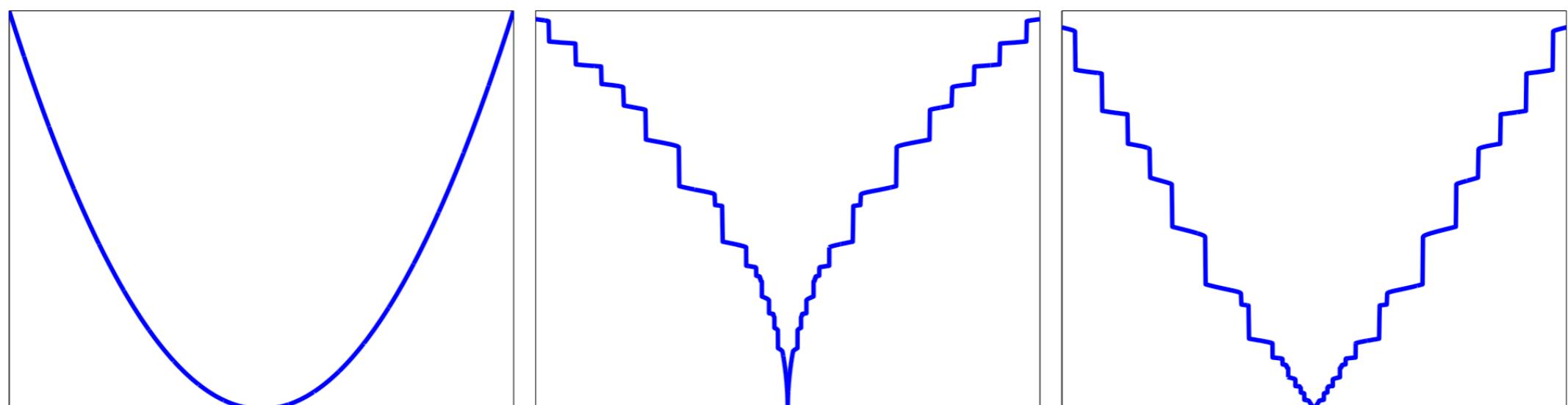
**Consequence:** same sequence  $\{\theta_t, t \geq 0\}$  on  $f$  or  $g \circ f$

# Comparison-based $\Rightarrow$ Invariance strict. increasing functions

Let  $g : \text{Im}f \rightarrow \text{Im}(g)$  be strictly increasing

$$f\left(X_{t+1}^{s(1)}\right) \leq \dots \leq f\left(X_{t+1}^{s(\lambda)}\right) \Leftrightarrow g \circ f\left(X_{t+1}^{s(1)}\right) \leq \dots \leq g \circ f\left(X_{t+1}^{s(\lambda)}\right)$$

**Consequence:** same sequence  $\{\theta_t, t \geq 0\}$  on  $f$  or  $g \circ f$



$$f(x) = \|x\|^2$$

$$g_1 \circ f$$

$$g_2 \circ f$$

CMA-ES - simplified setting  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$

## Sampling + ranking:

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

CMA-ES - simplified setting  $\theta_t = (m_t, \sigma_t, C_t) \in \mathbb{R}^n \times \mathbb{R}_> \times \mathcal{S}_{++}^n$

## Sampling + ranking:

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

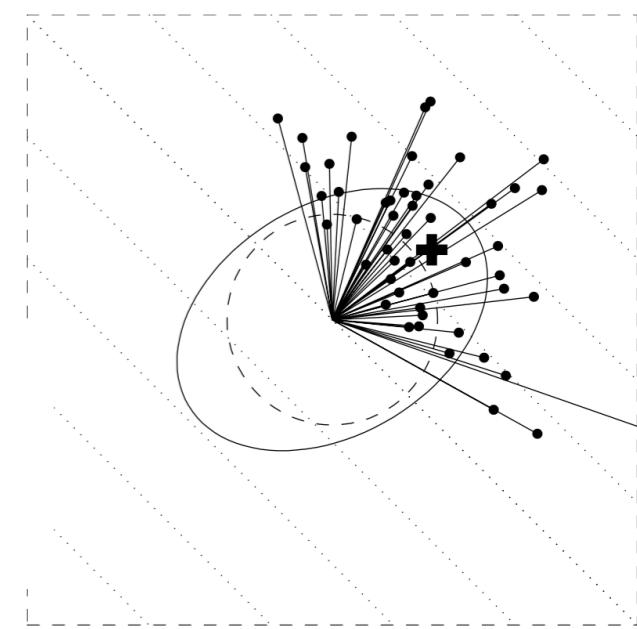
## Update of $\theta_t$ :

$$m_{t+1} = \sum_{i=1}^{\mu} w_i X_{t+1}^{s_{t+1}(i)} = m_t + \sigma_t \sqrt{C_t} \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)}$$

$$\sum_{i=1}^{\mu} w_i = 1, \mu_{\text{eff}} = 1 / \sum w_i^2$$

$$\sigma_{t+1} = \sigma_t \exp \left( \frac{c_\sigma}{d_\sigma} \left[ \frac{\sqrt{\mu_{\text{eff}}} \left\| \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)} \right\|}{E[\|\mathcal{N}(0, I_d)\|]} - 1 \right] \right)$$

$$C_{t+1} = (1 - c_\mu) C_t + c_\mu \sqrt{C_t} \underbrace{\left( \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)} [U_{t+1}^{s_{t+1}(i)}]^\top \right)}_{\text{rank } \mu \text{ update}} \sqrt{C_t}$$



# Cumulative Step-size Adaptation (CSA)

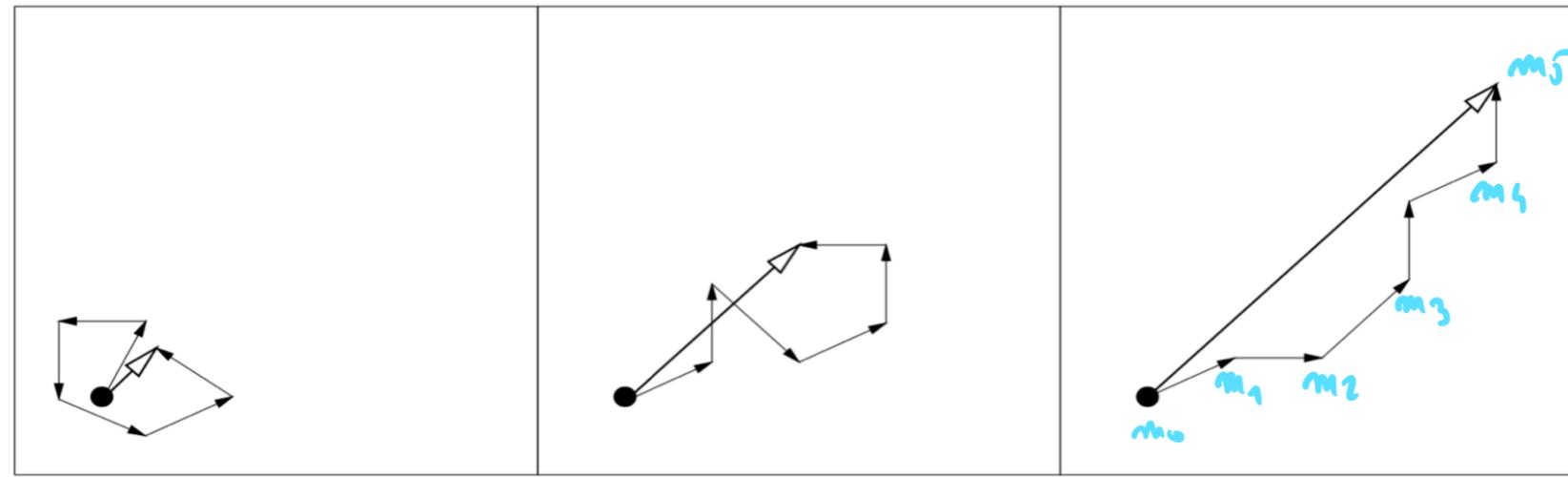
**Sampling + ranking:**

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i \quad i = 1, \dots, \lambda$$

$$\{U_t, t \geq 1\} \text{ i.i.d } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

**Idea:**



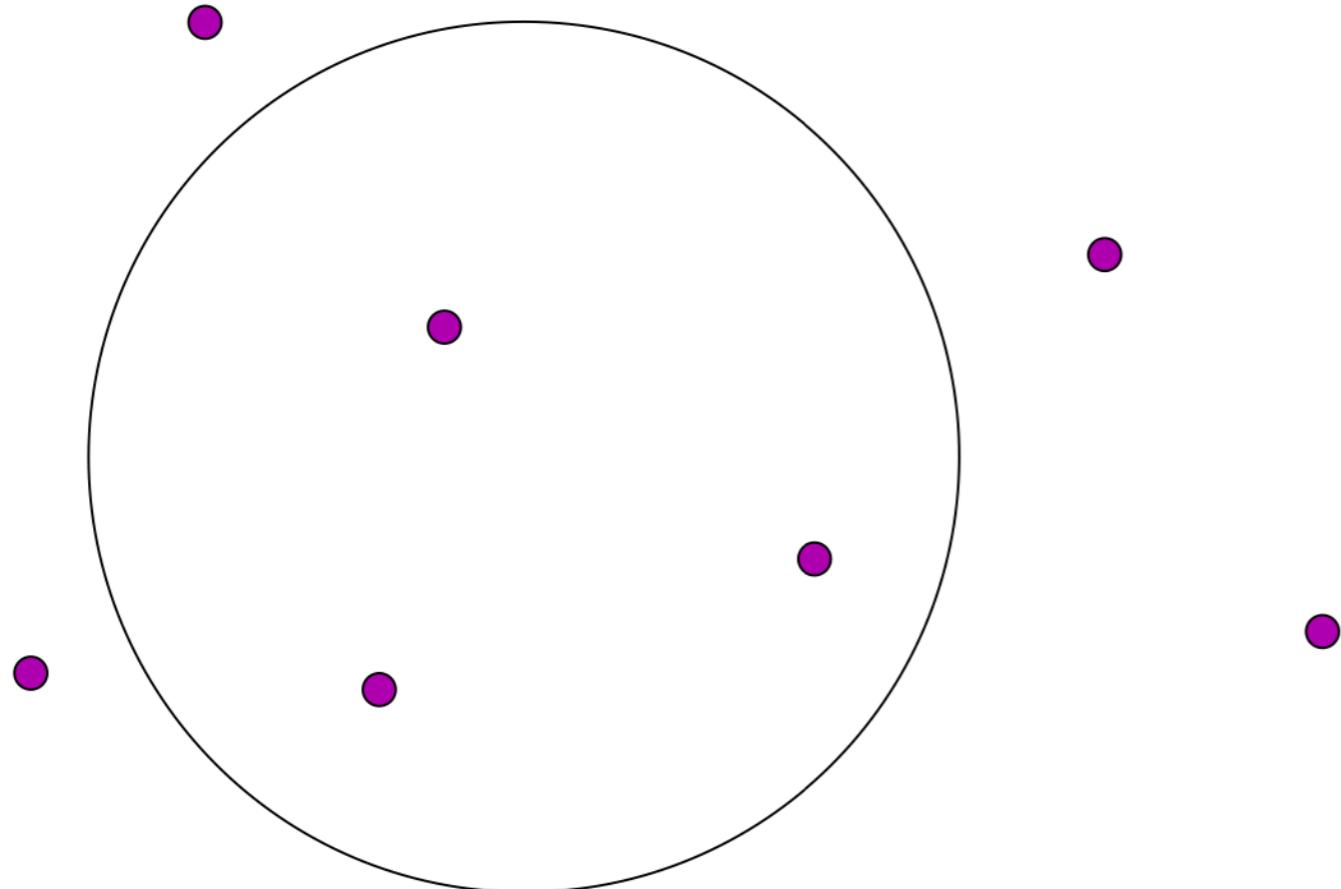
↓  
decrease  $\sigma$

↓  
increase  $\sigma$

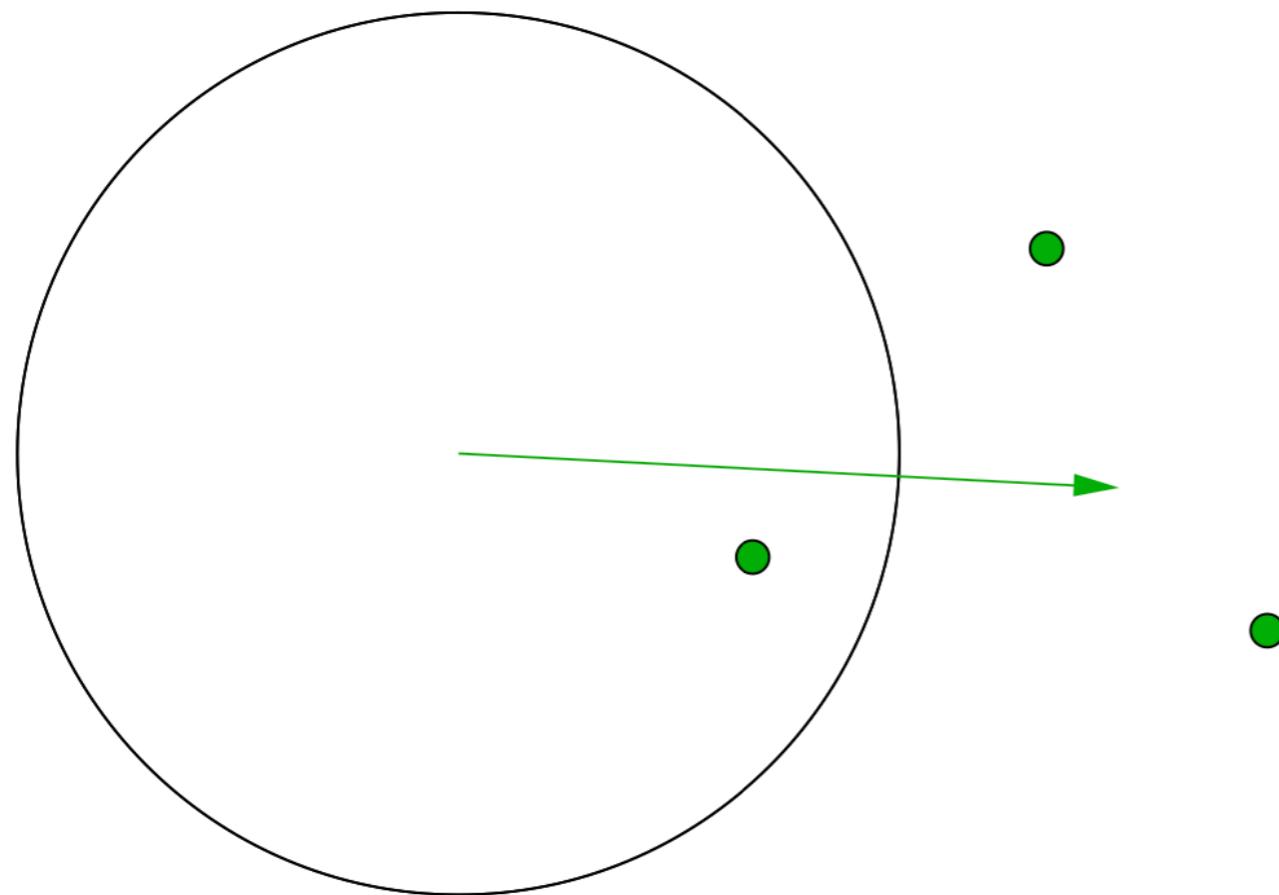
$$p_{t+1}^\sigma = (1 - c_\sigma) p_t^\sigma + \sqrt{c_\sigma(1 - c_\sigma)} \mu_{\text{eff}} \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)}$$

$$\sigma_{t+1} = \sigma_t \exp \left( \frac{c_\sigma}{d_\sigma} \left[ \frac{\|p_{t+1}^\sigma\|}{E[\|\mathcal{N}(0, I_d)\|]} - 1 \right] \right)$$

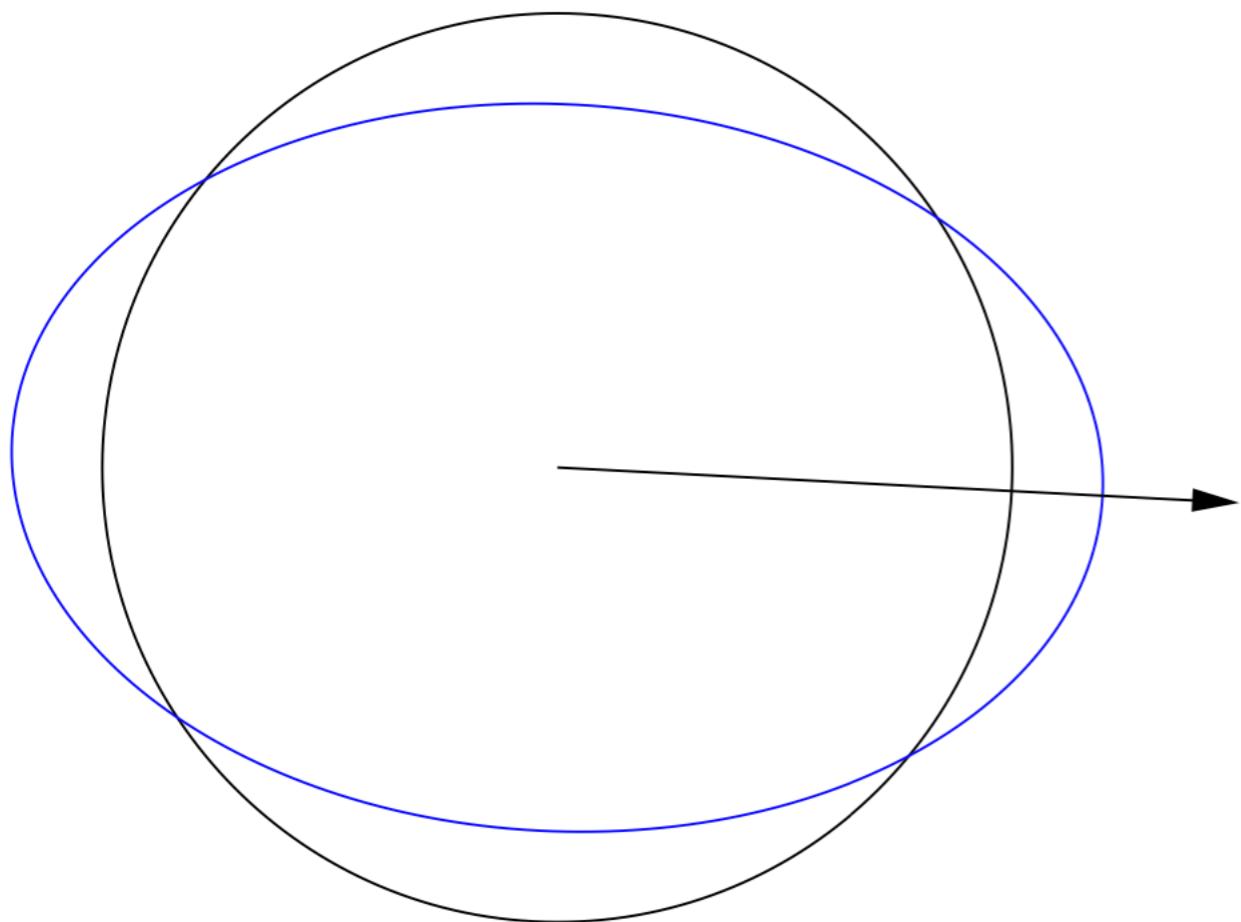
# Rank-one update

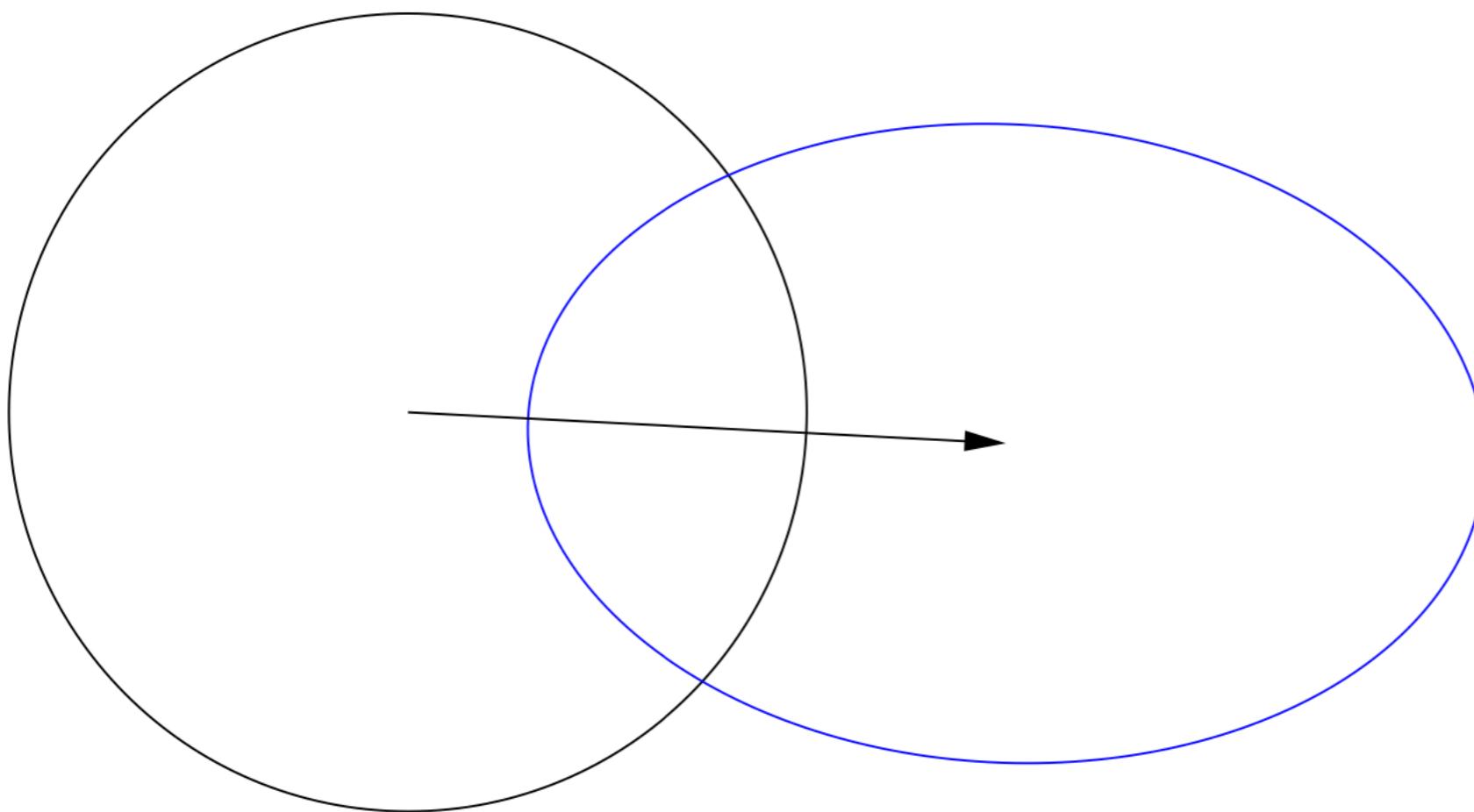


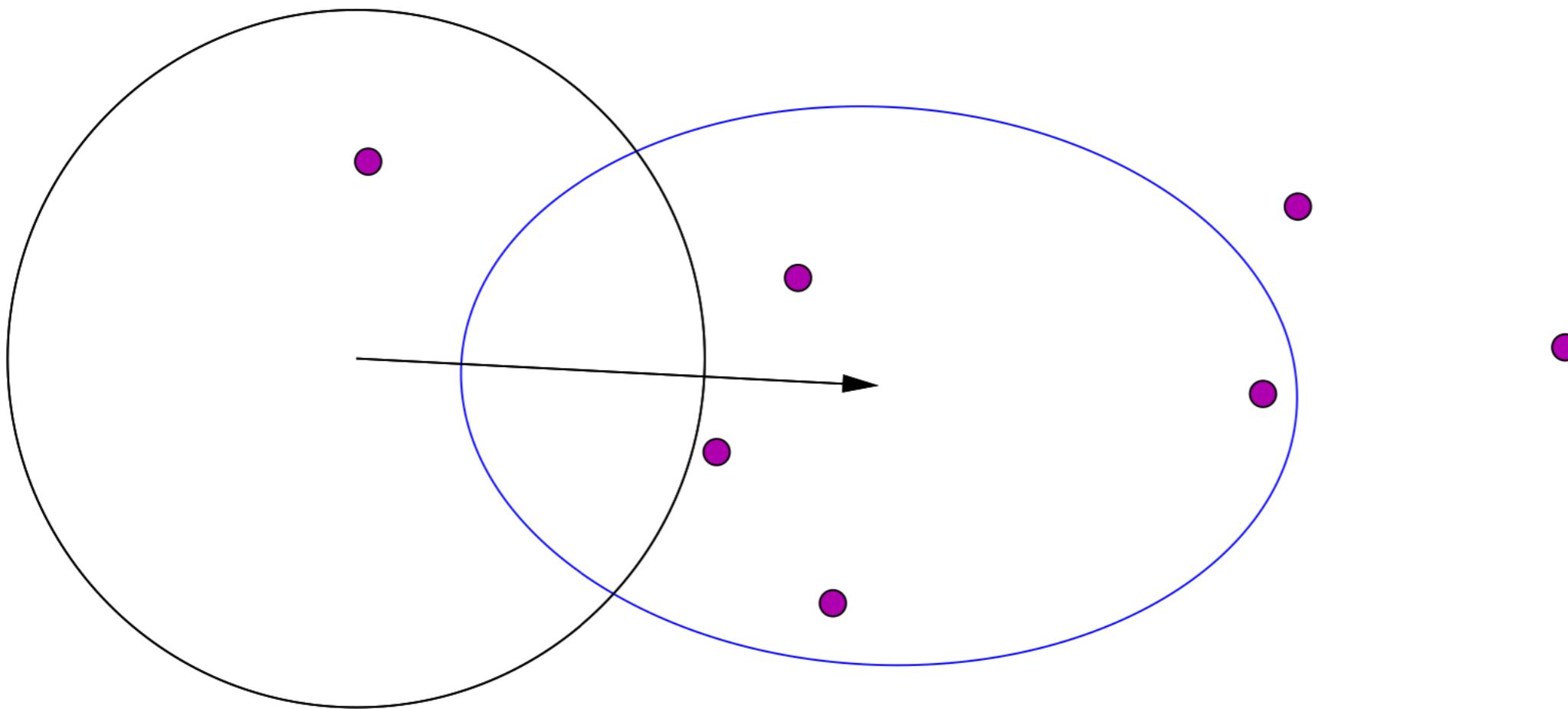
# Rank-one update

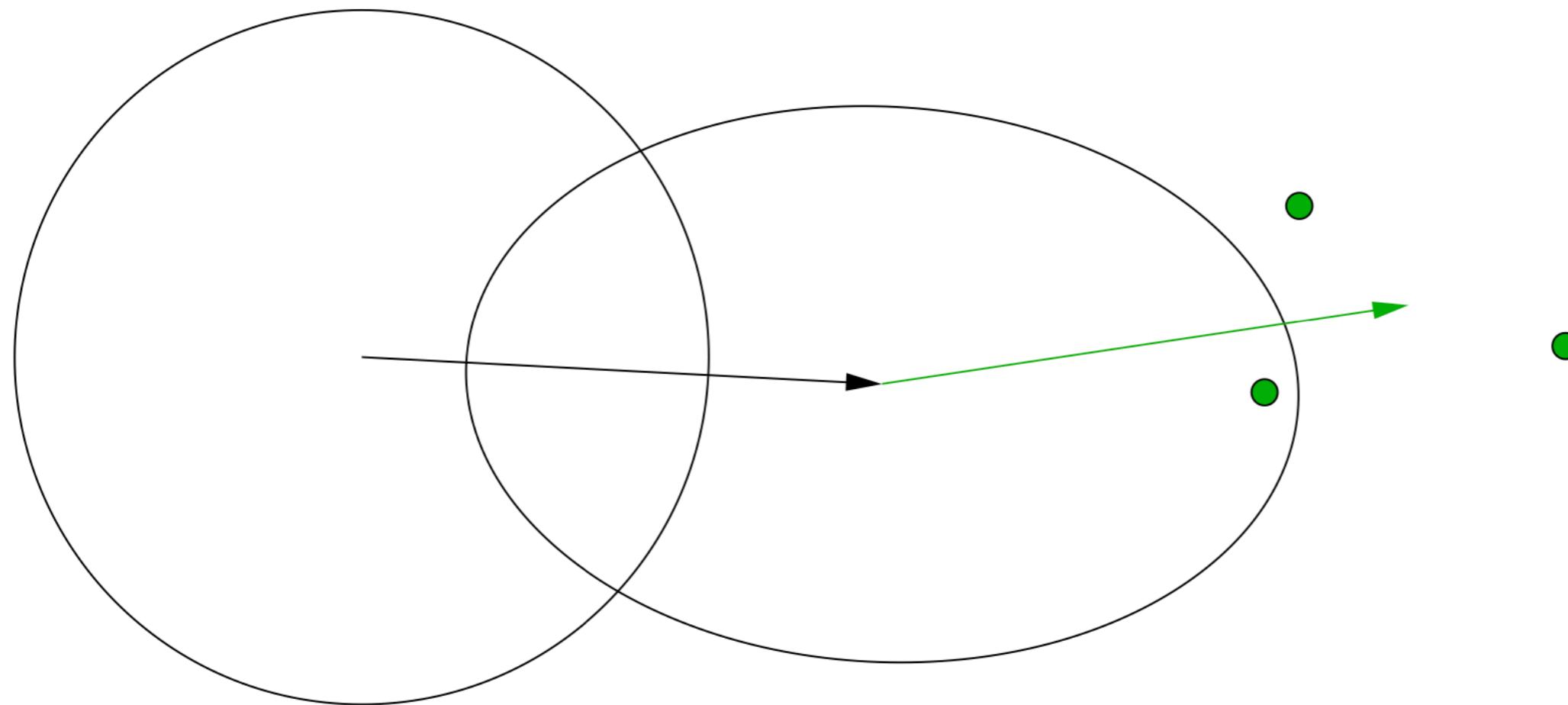


# Rank-one update









$$p_{t+1}^c = (1 - \textcolor{violet}{c}_c)p_t^c + \sqrt{\textcolor{violet}{c}_c(1 - \textcolor{violet}{c}_c)\mu_{\text{eff}}}\sqrt{C_t} \sum_{i=1}^{\mu} \textcolor{violet}{w}_i U_{t+1}^{s_{t+1}(i)}$$
$$C_{t+1} = (1 - \textcolor{violet}{c}_{\mu} - \textcolor{violet}{c}_1)C_t + \textcolor{violet}{c}_{\mu} \sqrt{C_t} \left( \sum_{i=1}^{\mu} \textcolor{violet}{w}_i U_{t+1}^{s_{t+1}(i)} [U_{t+1}^{s_{t+1}(i)}]^\top \right) \sqrt{C_t} + \textcolor{violet}{c}_1 \underbrace{\overbrace{p_{t+1}^c p_{t+1}^\top}^{\text{rank 1 update}}}$$

# CMA-ES

## Sampling + ranking:

$$X_{t+1}^i = m_t + \sigma_t \sqrt{C_t} U_{t+1}^i \quad i = 1, \dots, \lambda \quad \{U_t, t \geq 1\} \text{ i.i.d } U_{t+1}^i \sim \mathcal{N}(0, I_d)$$

$$f(X_{t+1}^{s_{t+1}(1)}) \leq \dots \leq f(X_{t+1}^{s_{t+1}(\lambda)})$$

## Update of $\theta_t$ :

$$m_{t+1} = \sum_{i=1}^{\mu} w_i X_{t+1}^{s_{t+1}(i)} = m_t + \sigma_t \sqrt{C_t} \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)}$$

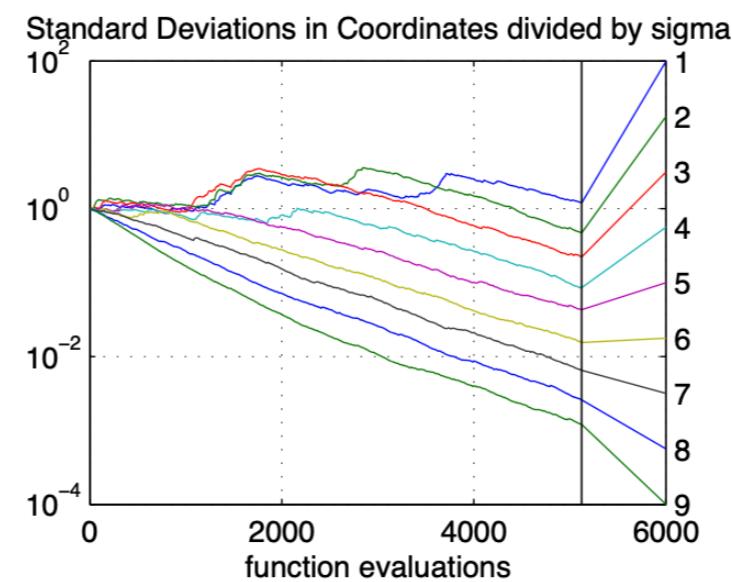
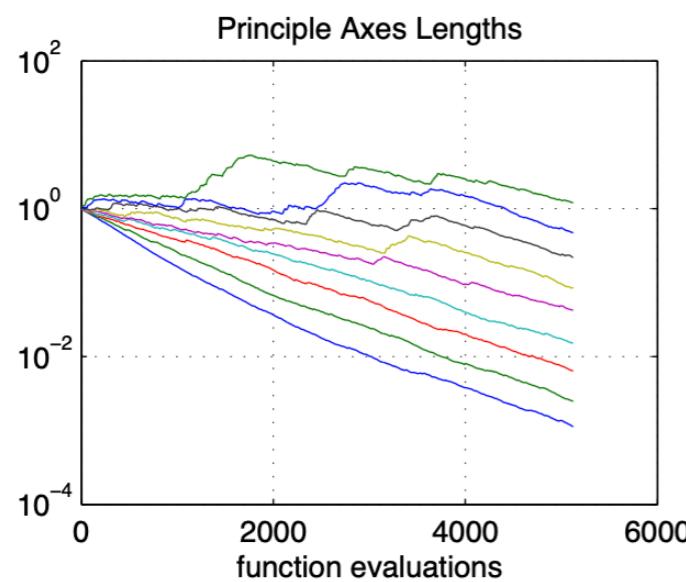
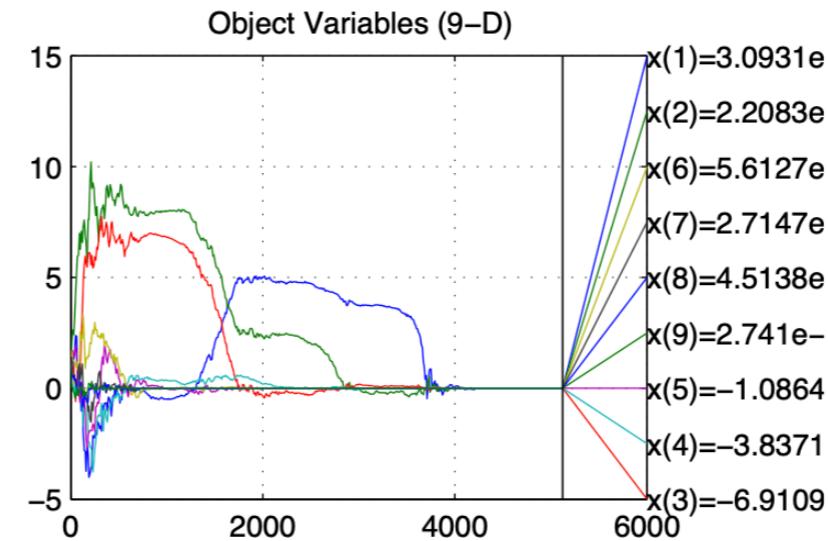
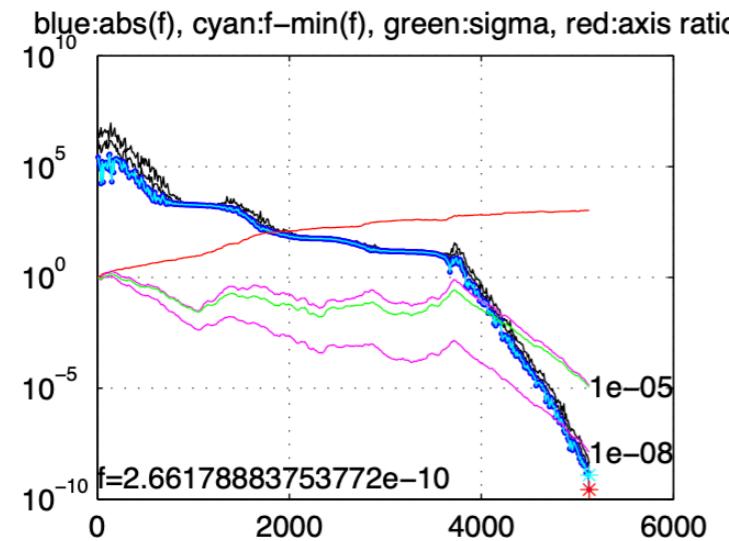
$$p_{t+1}^\sigma = (1 - c_\sigma) p_t^\sigma + \sqrt{c_\sigma(1 - c_\sigma)\mu_{\text{eff}}} \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)}$$

$$\sigma_{t+1} = \sigma_t \exp \left( \frac{c_\sigma}{d_\sigma} \left[ \frac{\|p_{t+1}^\sigma\|}{E[\|\mathcal{N}(0, I_d)\|]} - 1 \right] \right)$$

$$p_{t+1}^c = (1 - c_c) p_t^c + \sqrt{c_c(1 - c_c)\mu_{\text{eff}}} \sqrt{C_t} \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)}$$

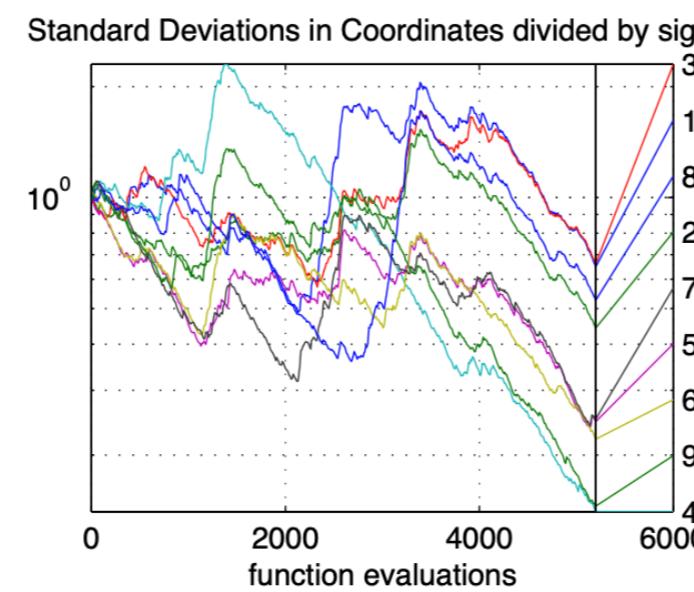
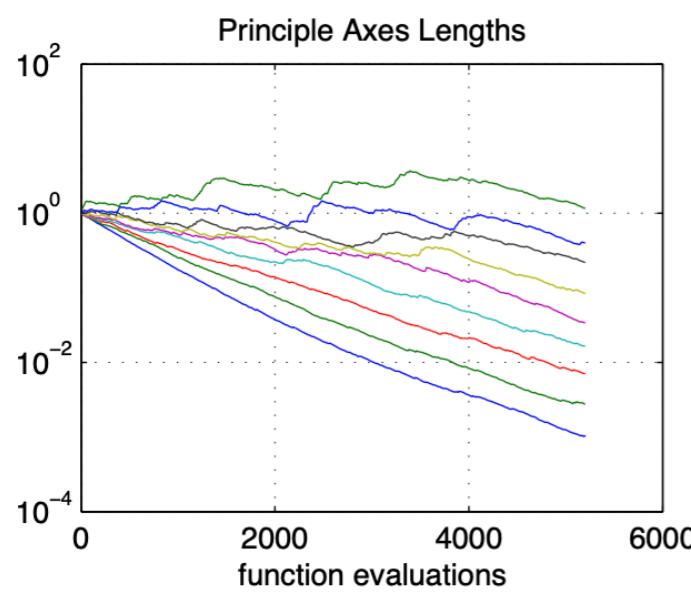
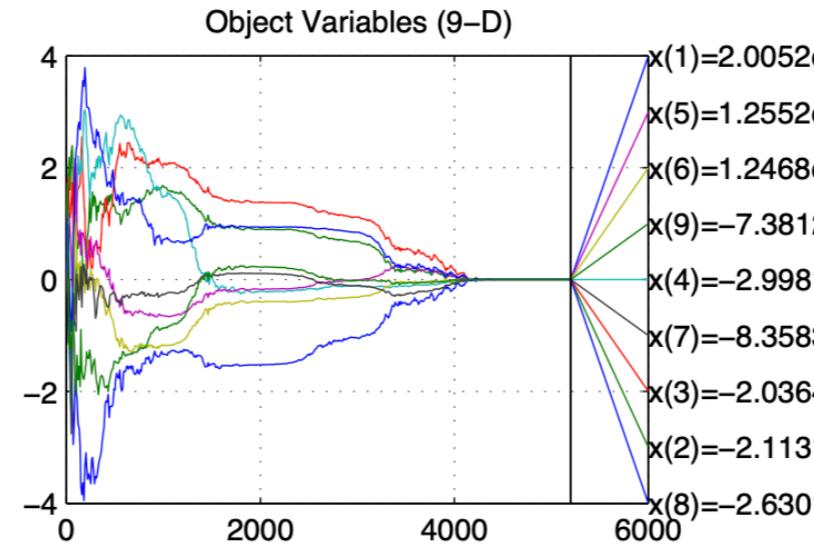
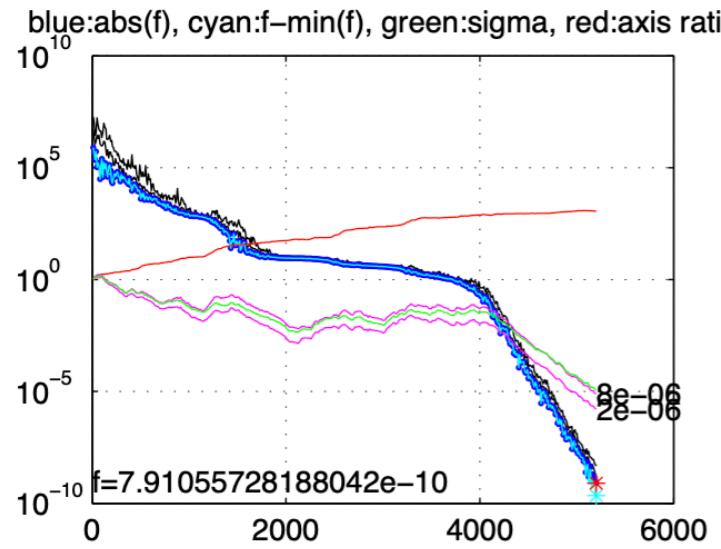
$$C_{t+1} = (1 - c_\mu - c_1) C_t + c_\mu \underbrace{\sqrt{C_t} \left( \sum_{i=1}^{\mu} w_i U_{t+1}^{s_{t+1}(i)} [U_{t+1}^{s_{t+1}(i)}]^\top \right) \sqrt{C_t}}_{\text{rank } \mu \text{ update}} + c_1 \underbrace{\overbrace{p_{t+1}^c p_{t+1}^\top}^{\text{rank 1 update}}}_{\text{rank 1 update}}$$

# Visualization - experimentum crucis



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

# Visualization - experimentum crucis



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

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

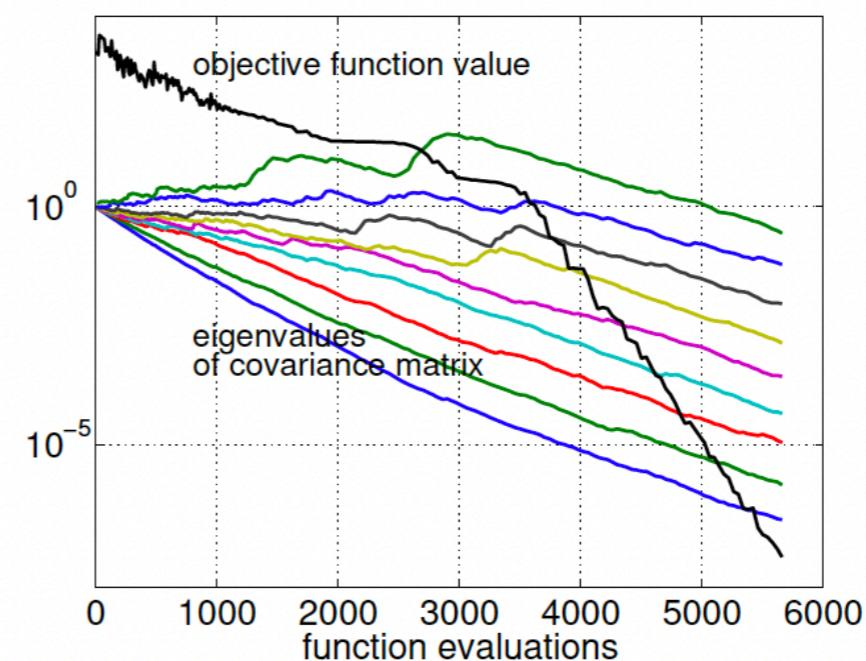
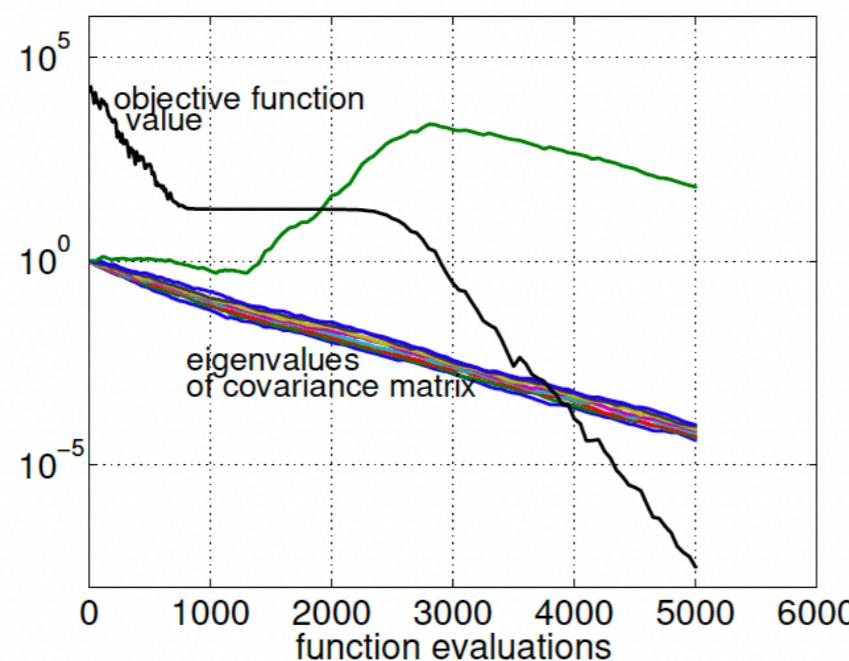
# Linear Convergence and Learning Inverse Hessian

For all  $g : \text{Im}(f) \rightarrow \mathbb{R}$ , strict increasing, for all  $f(x) = \frac{1}{2}(x - x^*)^\top H(x - x^*)$  with  $H \succ 0$  (SDP), when CMA-ES optimizes  $x \mapsto g \circ f(x)$ :

$$\frac{1}{t} \ln \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} \xrightarrow[t \rightarrow \infty]{} -\text{CR}$$

$$C_t \propto \alpha_t H^{-1} \text{ with } \alpha_t \rightarrow 0$$

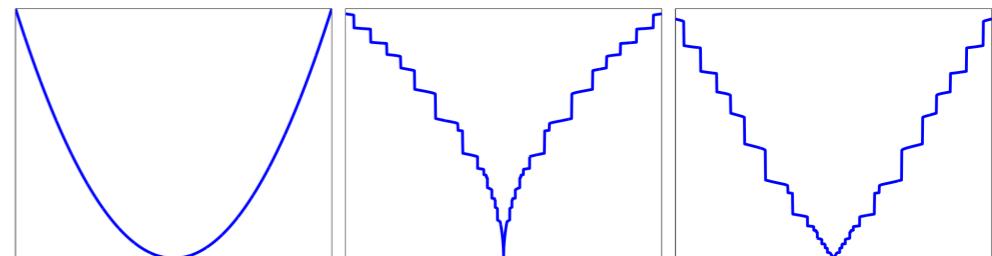
**Two examples:**



# Key of success (?) - Main features

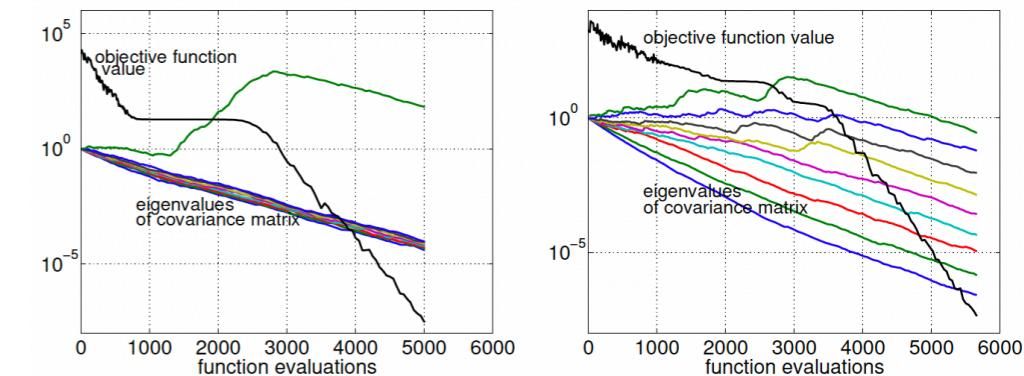
Invariances:

- to strictly increasing transformation of objective functions



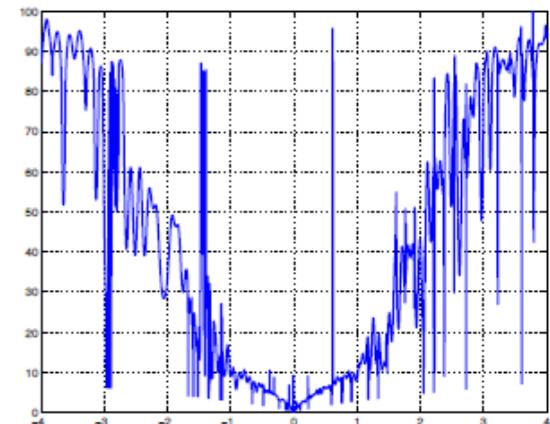
- affine invariance

$$x \mapsto \|x\|^2 \Leftrightarrow x \mapsto x^T H x \text{ for all } H \geq 0$$



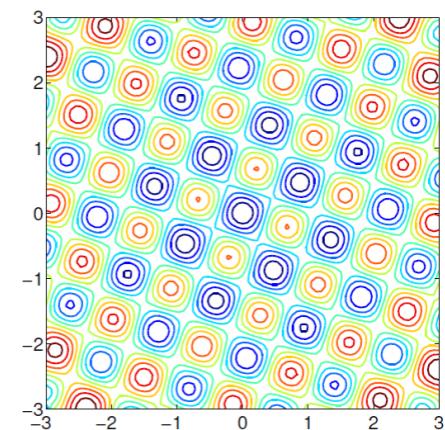
Control of diversity via step-size

- sees the function at the scale  $\sigma_t$   
**robust to noise, outliers, irregularities**



# Key of success (?) - Main features

Population-based algorithm: increase  $\lambda$  makes algorithm more global



Careful tuning of all constants (depending on dimension, ...)  
**parameter-free algorithm**

Original designers did not care about having a convergence proof

- **successful approach**: framework not restricted to updates for which proofs can be made
- was challenging to prove linear convergence of CMA-ES [PhD thesis Armand Gissler, 2024]

Thanks!