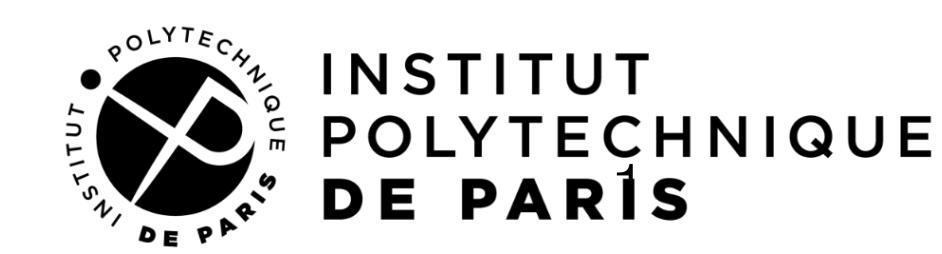


It's Always the Step-Size

A report of a semi-practitioner

Nikolaus Hansen
Inria & École polytechnique, France

February 2025



Overview

- empirical research (my methodology)
- what is cGAs K in CMA-ES?
 K is considered to be a replacement/analog for population size, $1/K$ is a learning rate
- in Evolution Strategies (ES)
 - the optimal standard deviation for sampling is proportional to the parent number μ
 - diversity is close to divergence
 - premature diversity loss due to genetic drift happens in noisy landscapes
 - it's always the step-size

The CMA-ES is an EDA

Input: $\mathbf{m} \in \mathbb{R}^n$; $\sigma \in \mathbb{R}_+$; $\lambda \in \mathbb{N}_{\geq 2}$, usually $\lambda \geq 5$, default $4 + \lfloor 3 \log n \rfloor$

Initialize $\mathbf{C} = \mathbf{I}$

While not *terminate*

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

$$\mathbf{m} \leftarrow \mathbf{m} + c_m \sigma \mathbf{y}_w, \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} w_{\text{rk}(i)} \mathbf{y}_i \quad \text{update mean}$$

$$\mathbf{C} \leftarrow \mathbf{C} + c_\mu \sum_{i=1}^{\lambda} w_{\text{rk}(i)} (\mathbf{y}_i \mathbf{y}_i^\top - \mathbf{C}) \quad \text{update C}$$

The CMA-ES is an EDA

Input: $\mathbf{m} \in \mathbb{R}^n$; $\sigma \in \mathbb{R}_+$; $\lambda \in \mathbb{N}_{\geq 2}$, usually $\lambda \geq 5$, default $4 + \lfloor 3 \log n \rfloor$

Set $c_m = 1$; $c_1 \approx 2/n^2$; $c_\mu \approx \mu_w/n^2$; $c_c \approx 4/n$; $c_\sigma \approx 1/\sqrt{n}$; $d_\sigma \approx 1$; $w_{i=1\dots\lambda}$ decreasing in i and $\sum_i^\mu w_i = 1$, $w_\mu > 0 \geq w_{\mu+1}$, $\mu_w^{-1} := \sum_{i=1}^\mu w_i^2 \approx 3/\lambda$

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

While not *terminate*

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

$$\mathbf{m} \leftarrow \mathbf{m} + c_m \sigma \mathbf{y}_w, \quad \text{where } \mathbf{y}_w = \sum_{i=1}^\mu w_{\text{rk}(i)} \mathbf{y}_i \quad \text{update mean}$$

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

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbf{1}_{[0,2n]} \left\{ \|\mathbf{p}_\sigma\|^2 \right\} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w \quad \text{path for } \mathbf{C}$$

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

$$\mathbf{C} \leftarrow \mathbf{C} + c_\mu \sum_{i=1}^\lambda w_{\text{rk}(i)} (\mathbf{y}_i \mathbf{y}_i^\top - \mathbf{C}) + c_1 (\mathbf{p}_c \mathbf{p}_c^\top - \mathbf{C}) \quad \text{update } \mathbf{C}$$

Not covered: termination, restarts, useful output, search boundaries and encoding, corrections for: positive definiteness guaranty, \mathbf{p}_c variance loss, c_σ and d_σ for large λ

The CMA-ES

Input: $\mathbf{m} \in \mathbb{R}^n$; $\sigma \in \mathbb{R}_+$; $\lambda \in \mathbb{N}_{\geq 2}$, usually $\lambda \geq 5$, default $4 + \lfloor 3 \log n \rfloor$

Set $c_m = 1$; $c_1 \approx 2/n^2$; $c_\mu \approx \mu_w/n^2$; $c_c \approx 4/n$; $c_\sigma \approx 1/\sqrt{n}$; $d_\sigma \approx 1$; $w_{i=1\dots\lambda}$ decreasing in i and $\sum_i^\mu w_i = 1$, $w_\mu > 0 \geq w_{\mu+1}$, $\mu_w^{-1} := \sum_{i=1}^\mu w_i^2 \approx 3/\lambda$

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

c_m, c_μ, c_1 are analog
to $1/K$ in cGA,
 d_σ is analog to K

While not *terminate*

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

$$\mathbf{m} \leftarrow \mathbf{m} + \underbrace{c_m}_{} \sigma \mathbf{y}_w, \quad \text{where } \mathbf{y}_w = \sum_{i=1}^\mu w_{\text{rk}(i)} \mathbf{y}_i \quad \text{update mean}$$

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

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbf{1}_{[0,2n]} \left\{ \|\mathbf{p}_\sigma\|^2 \right\} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w \quad \text{path for } \mathbf{C}$$

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

$$\mathbf{C} \leftarrow \mathbf{C} + \underbrace{c_\mu}_{} \sum_{i=1}^\lambda w_{\text{rk}(i)} (\mathbf{y}_i \mathbf{y}_i^\top - \mathbf{C}) + \underbrace{c_1}_{} (\mathbf{p}_c \mathbf{p}_c^\top - \mathbf{C}) \quad \text{update } \mathbf{C}$$

Not covered: termination, restarts, useful output, search boundaries and encoding, corrections for: positive definiteness guaranty, \mathbf{p}_c variance loss, c_σ and d_σ for large λ

The CMA-ES

Input: $\mathbf{m} \in \mathbb{R}^n$; $\sigma \in \mathbb{R}_+$; $\lambda \in \mathbb{N}_{\geq 2}$, usually $\lambda \geq 5$, default $4 + \lfloor 3 \log n \rfloor$

Set $c_m = 1$; $c_1 \approx 2/n^2$; $c_\mu \approx \mu_w/n^2$; $c_c \approx 4/n$; $c_\sigma \approx 1/\sqrt{n}$; $d_\sigma \approx 1$; $w_{i=1\dots\lambda}$ decreasing in i and $\sum_i^\mu w_i = 1$, $w_\mu > 0 \geq w_{\mu+1}$, $\mu_w^{-1} := \sum_{i=1}^\mu w_i^2 \approx 3/\lambda$

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

While not terminate

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

$$\mathbf{m} \leftarrow \mathbf{m} + c_m \sigma \mathbf{y}_w, \quad \text{where } \mathbf{y}_w = \sum_{i=1}^\mu w_{\text{rk}(i)} \mathbf{y}_i \quad \text{update mean}$$

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

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbf{1}_{[0,2n]} \left\{ \|\mathbf{p}_\sigma\|^2 \right\} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w \quad \text{path for } \mathbf{C}$$

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

$$\mathbf{C} \leftarrow \mathbf{C} + c_\mu \sum_{i=1}^\lambda w_{\text{rk}(i)} (\mathbf{y}_i \mathbf{y}_i^\top - \mathbf{C}) + c_1 (\mathbf{p}_c \mathbf{p}_c^\top - \mathbf{C}) \quad \text{update } \mathbf{C}$$

c_m, c_μ, c_1 are analog
to $1/K$ in cGA,
 d_σ is analog to K

Not covered: termination, restarts, useful output, search boundaries and encoding,
corrections for: positive definiteness guaranty, \mathbf{p}_c variance loss, c_σ and d_σ for large λ

The CMA-ES

Input: $\mathbf{m} \in \mathbb{R}^n$; $\sigma \in \mathbb{R}_+$; $\lambda \in \mathbb{N}_{\geq 2}$, usually $\lambda \geq 5$, default $4 + \lfloor 3 \log n \rfloor$

Set $c_m = 1$; $c_1 \approx 2/n^2$; $c_\mu \approx \mu_w/n^2$; $c_c \approx 4/n$; $c_\sigma \approx 1/\sqrt{n}$; $d_\sigma \approx 1$; $w_{i=1\dots\lambda}$ decreasing in i and $\sum_i^\mu w_i = 1$, $w_\mu > 0 \geq w_{\mu+1}$, $\mu_w^{-1} := \sum_{i=1}^\mu w_i^2 \approx 3/\lambda$

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

c_m, c_μ, c_1 are analog
to $1/K$ in cGA,
 d_σ is analog to K

While not terminate

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

$$\mathbf{m} \leftarrow \mathbf{m} + c_m \sigma \mathbf{y}_w, \quad \text{where } \mathbf{y}_w = \sum_{i=1}^\mu w_{\text{rk}(i)} \mathbf{y}_i \quad \text{update mean}$$

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

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbf{1}_{[0,2n]} \left\{ \|\mathbf{p}_\sigma\|^2 \right\} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w \quad \text{path for } \mathbf{C}$$

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

$$\mathbf{C} \leftarrow \mathbf{C} + c_\mu \sum_{i=1}^\lambda w_{\text{rk}(i)} (\mathbf{y}_i \mathbf{y}_i^\top - \mathbf{C}) + c_1 (\mathbf{p}_c \mathbf{p}_c^\top - \mathbf{C}) \quad \text{update } \mathbf{C}$$

Not covered: termination, restarts, useful output, search boundaries and encoding, corrections for: positive definiteness guaranty, \mathbf{p}_c variance loss, c_σ and d_σ for large λ

```

1 %matplotlib widget
2 import cma
3
4 x, es = cma.fmin2(cma.ff.rosen, 14 * [1], 1, {'CMA_diagonal_decoding': True, 'ftarget': 1e-7})
5 es.plot();

```

✓ 2.2s

Python

(5_w,11)-aCMA-ES ($\mu_w=3.4, w_1=42\%$) in dimension 14 (seed=553091, Mon Feb 24 13:19:30 2025)

Iterat #Fevals function value axis ratio sigma min&max std t[m:s]

1	11	1.040197601953083e+03	1.0e+00	8.93e-01	7e-01	9e-01	0:00.0
2	22	1.476140839356284e+03	1.1e+00	8.01e-01	6e-01	8e-01	0:00.0
3	33	1.027303494859089e+03	1.1e+00	7.42e-01	6e-01	8e-01	0:00.0
100	1100	1.331908726689222e-01	3.7e+00	3.99e-02	4e-03	1e-02	0:00.1
200	2200	7.625477142936609e-04	2.4e+01	7.53e-03	3e-04	8e-03	0:00.2
265	2915	9.436810372489095e-08	3.0e+01	1.86e-04	6e-06	2e-04	0:00.2

termination on ftarget=1e-07 (Mon Feb 24 13:19:31 2025)

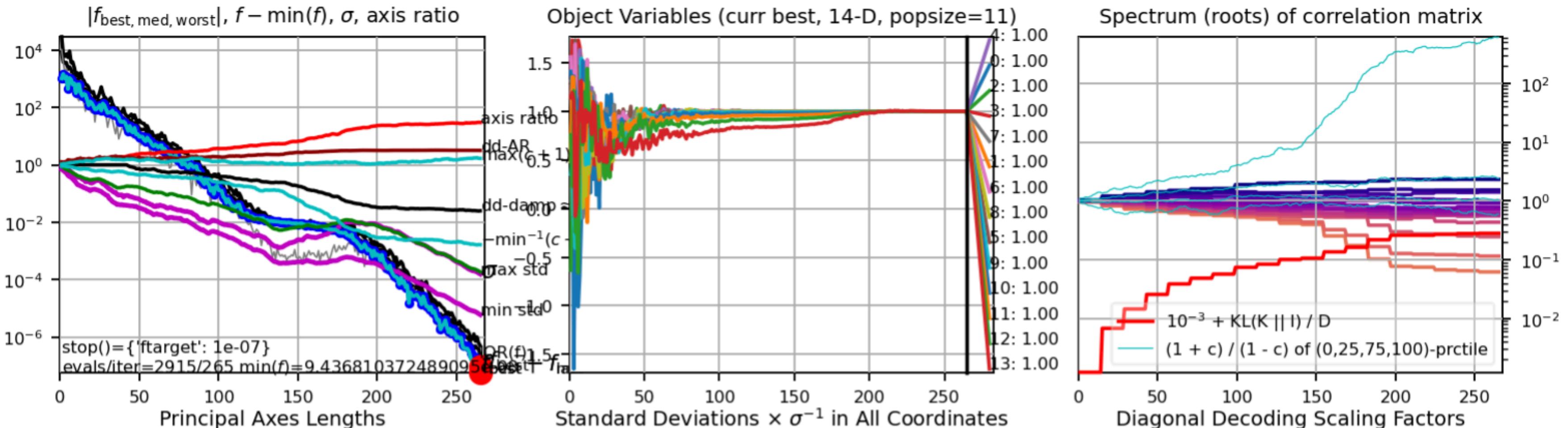
final/bestever f-value = 3.152852e-08 3.152852e-08 after 2916/2916 evaluations

incumbent solution: [1.00000409 1.00000359 1.0000048 1.00000458 1.00000313 1.00000012
0.99999997 1.00000003 ...]

std deviations: [6.00652667e-06 6.04096130e-06 6.75724117e-06 6.71399653e-06

6.68393736e-06 6.42605509e-06 6.86269378e-06 6.56077278e-06 ...]

Figure 325



```

3 55 1.027363494059009e+03 1.1e+00 7.42e-01 0e-01 0e-01 0.0000
100 1100 1.331908726689222e-01 3.7e+00 3.99e-02 4e-03 1e-02 0:00.1
200 2200 7.625477142936609e-04 2.4e+01 7.53e-03 3e-04 8e-03 0:00.2
265 2915 9.436810372489095e-08 3.0e+01 1.86e-04 6e-06 2e-04 0:00.2

```

termination on ftarget=1e-07 (Mon Feb 24 13:19:31 2025)

final/bestever f-value = 3.152852e-08 3.152852e-08 after 2916/2916 evaluations

incumbent solution: [1.00000409 1.00000359 1.0000048 1.00000458 1.00000313 1.00000012

0.99999997 1.00000003 ...]

std deviations: [6.00652667e-06 6.04096130e-06 6.75724117e-06 6.71399653e-06

6.68393736e-06 6.42605509e-06 6.86269378e-06 6.56077278e-06 ...]

Figure 325

