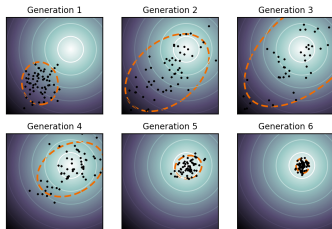


Optimization in Machine Learning

CMA-ES Algorithm



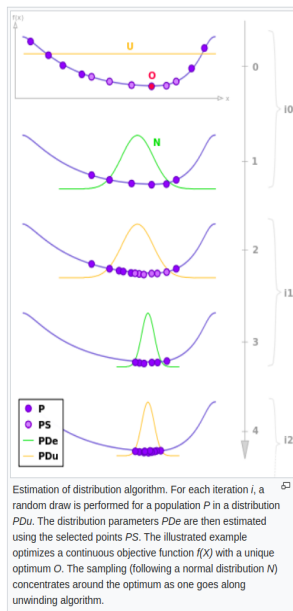
Learning goals

- CMA-ES strategy
- Estimation of distribution
- Step size control

CMA-ES

General algorithmic template with initial parameter setting $\theta^{[0]}$ of parameterized density $p(\mathbf{x}|\theta)$, such that each iteration $t \in \{0, 1, \dots, T\}$ consist of:

- 1 Draw λ samples, $\mathbf{x}^{(k)}$ from $p(\mathbf{x}|\theta^{[t]})$
- 2 Evaluate $W(f(\mathbf{x}^{(k)}))$, where $W(\cdot)$ gives weights for each $\mathbf{x}^{(k)}$, typically 0 or 1 (order-preserving fitness transformation)
- 3 Find a $\theta^{[t+1]}$ that uses weighted samples and corresponding function evaluations to move $p(\mathbf{x}|\theta)$ towards regions \mathcal{S} that have large function values.



CMA-ES

Covariance Matrix Adaptation Evolution Strategy (CMA-ES) is

- A state-of-the-art tool in evolutionary computation
- A stochastic/randomized method
- For usage in continuous domain
- For non-linear, non-convex optimization problems
- Useful in case “classical” search methods like quasi-Newton methods (BFGS) or conjugate gradient methods fail due to a non-convex or rugged search landscape (e.g. outliers, noise, local optima, sharp bends).

Detailed information on CMA-ES can be found in

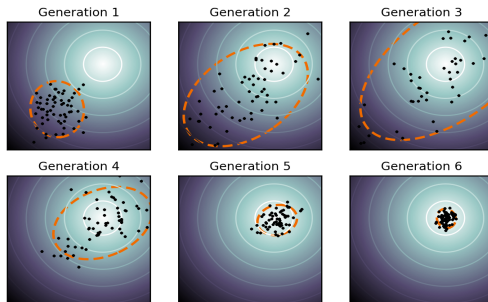
- ❶ Nikolaus Hansen. The CMA Evolution Strategy. 2005
- ❷ A. Auger, N. Hansen: Tutorial CMA-ES: Evolution Strategies and Covariance Matrix Adaptation. 2012.

CMA-ES: STRATEGY

A population of new search points (individuals, offspring) is generated by sampling a multivariate normal distribution. A repeated update of the mean vector and covariance matrix with the respectively best ranked individuals moves the distributions towards the optimum.

Search points for generation number $t = 0, 1, \dots, T$:

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



CMA-ES: STRATEGY

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

- $\mathbf{x}^{[t+1](k)} \in \mathbb{R}^d$ is the k -th offspring (search point, individual) from generation $t + 1$.
- $\lambda \geq 2$ is the population/sample size, number of offsprings.

At generation t :

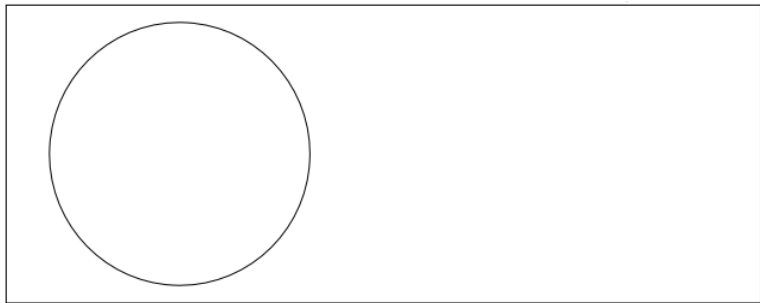
- $\mathcal{N}(\mathbf{0}, \mathbf{C}^{[t]})$ is multivariate normal distribution with zero mean, covariance matrix $\mathbf{C}^{[t]}$. *Note:*
 $\mathbf{m}^{[t]} + \sigma^{[t]} \mathcal{N}(\mathbf{0}, \mathbf{C}^{[t]}) \sim \mathcal{N}(\mathbf{m}^{[t]}, (\sigma^{[t]})^2 \mathbf{C}^{[t]}).$
- $\mathbf{m}^{[t]} \in \mathbb{R}^d$ is the mean value of the search distribution.
- $\sigma^{[t]} \in \mathbb{R}_+$ is the “overall” standard deviation/step size.
- $\mathbf{C}^{[t]} \in \mathbb{R}^{d \times d}$ is the covariance matrix.

→ *How to calculate $\mathbf{m}^{[t+1]}$, $\mathbf{C}^{[t+1]}$, $\sigma^{[t+1]}$ for next generation $t + 1$?*

CMA-ES: BASIC METHOD - ITERATION 1

❶ Sample from distribution

$\mathbf{x}^{[1](k)} = \mathbf{m}^{[0]} + \sigma^{[0]} \mathcal{N}(\mathbf{0}, \mathbf{C}^{[0]})$ multivariate normal distribution.

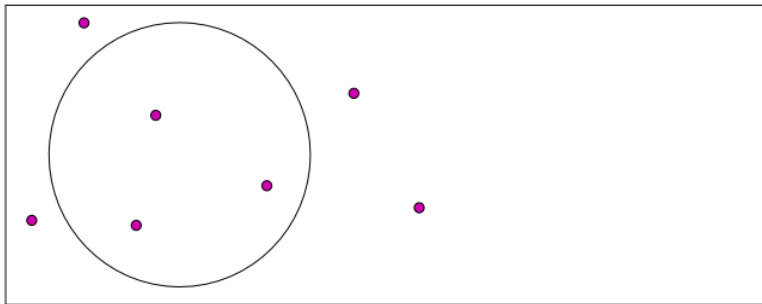


Initial distribution $\mathcal{N}^{[0]} \sim (\mathbf{0}, \mathbb{I}_2)$ of generation $t = 0$.

CMA-ES: BASIC METHOD - ITERATION 1

1 Sample from distribution

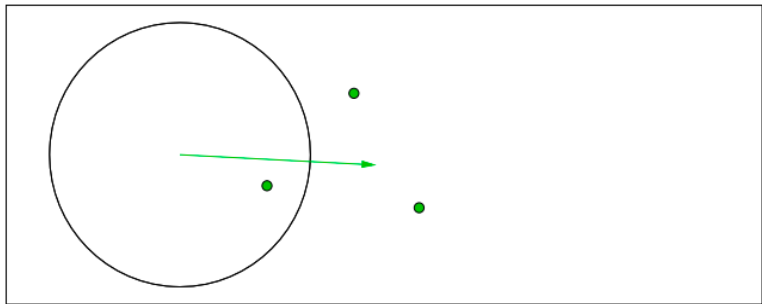
$\mathbf{x}^{[1](k)} = \mathbf{m}^{[0]} + \sigma^{[0]} \mathcal{N}(\mathbf{0}, \mathbf{C}^{[0]})$ multivariate normal distribution.



Initial distribution $\mathcal{N}^{[0]} \sim (\mathbf{0}, \mathbb{I}_2)$ of generation $t = 0$, $\lambda = 6$.

CMA-ES: BASIC METHOD - ITERATION 1

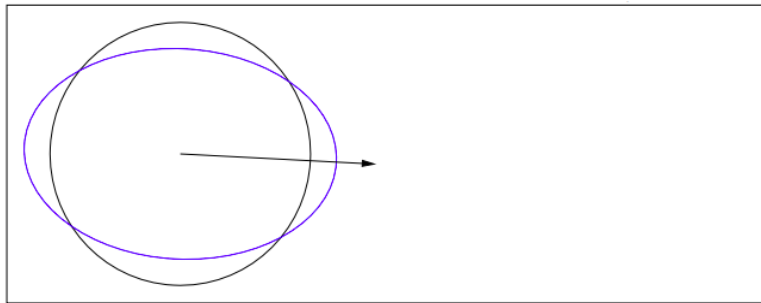
- ② **Ranking** solutions according to their fitness (*Selection* of μ best)
 $\mathbf{x}_{i:\lambda}$ as i -th ranked solution point, such that $f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.



Calculation of auxiliary variable $\mathbf{y}_w := \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$, using $\mu = 3$
selected points (high fitness \rightarrow high weights)

CMA-ES: BASIC METHOD - ITERATION 1

- ③ **Update covariance matrix** (*Recombination*),
improving “expected fitness” and likelihood for good steps.

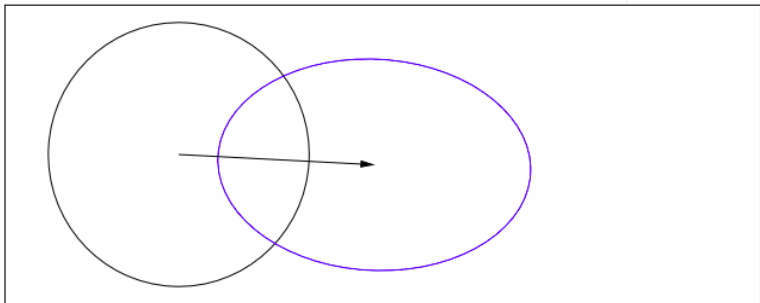


Blue circle as a mixture of \mathbf{C} and step \mathbf{y}_w (simplified):

$$\mathbf{C}^{[1]} = 0.8\mathbf{C}^{[0]} + 0.2\mathbf{y}_w^{[0]}(\mathbf{y}_w^{[0]})^\top \text{ (Rank 1 update).}$$

CMA-ES: BASIC METHOD - ITERATION 1

4 Update mean

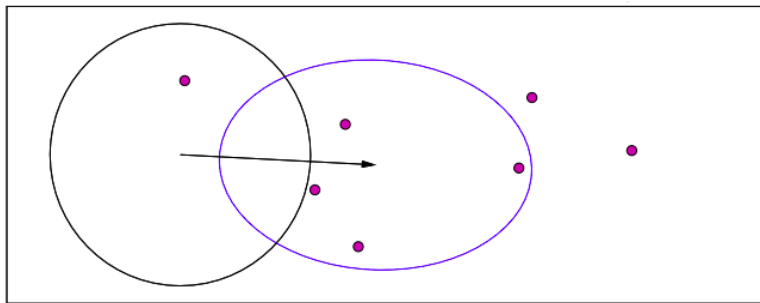


Movement towards the new distribution with mean

$$\mathbf{m}^{[1]} = \mathbf{m}^{[0]} + \sigma^{[0]} \mathbf{y}_w^{[0]}.$$

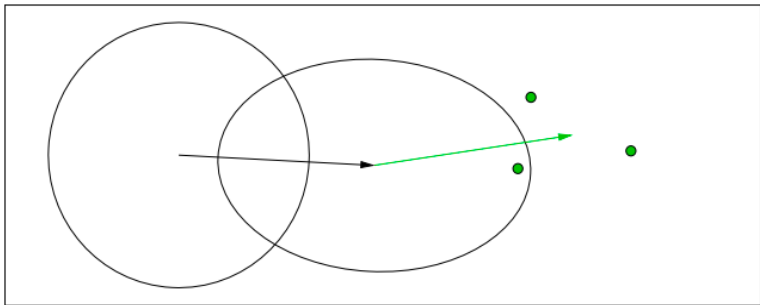
CMA-ES: BASIC METHOD - ITERATION 2

- 1 **Sample** from distribution for new generation



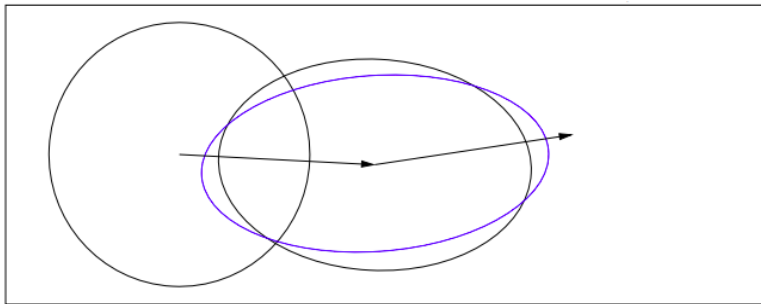
CMA-ES: BASIC METHOD - ITERATION 2

- ② **Ranking** solutions according to their fitness (*Selection* of μ best)



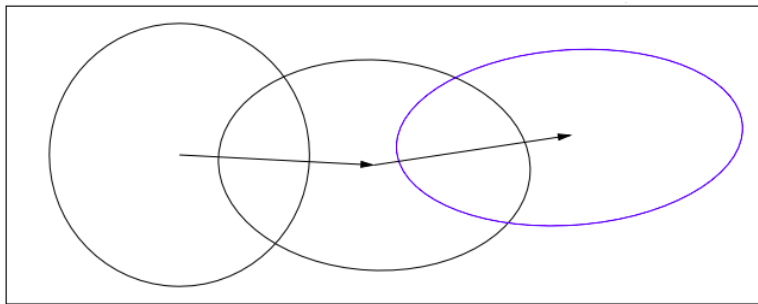
CMA-ES: BASIC METHOD - ITERATION 2

③ Update mean and covariance matrix (*Recombination*)



CMA-ES: BASIC METHOD - ITERATION 2

- ④ **Update step-size** based on non-local information, exploit correlations in the history of steps.



UPDATING \mathbf{C} : CMA - RANK-ONE UPDATE

Initialize $\mathbf{m} \in \mathbb{R}^d$ and $\mathbf{C} = \mathbb{I}$, set $\sigma = 1$, learning rate $c_{cov} \approx 2/d^2$.
While not terminate

$$\mathbf{x}^{(k)} = \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C})$$

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

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

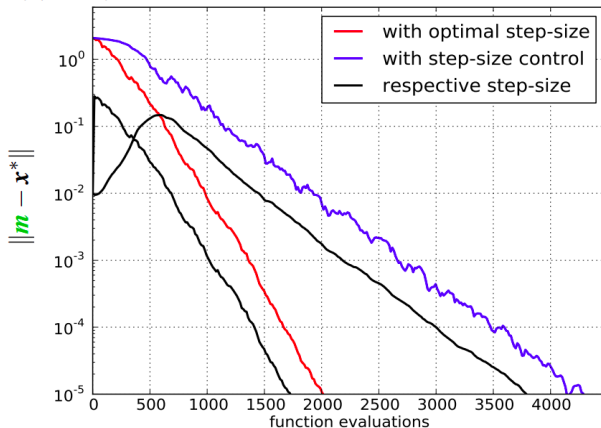
The rank-one update was developed in several domains independently, conducting a **principle component analysis** (PCA) of steps \mathbf{y}_w sequentially in time and space.

UPDATING σ : METHODS STEP-SIZE CONTROL

- **1/5-th success rule**: increases the step-size if more than 20 % of the new solutions are successful, decrease otherwise
- **σ -self-adaptation**: mutation is applied to the step-size and the better - according to the objective function value - is selected
- **Path length control via cumulative step-size adaptation (CSA)**: self-adaptation derandomized and non-localized
- Alternative step-size adaptation mechanism: two-point step-size adaptation, median success rule, population success rule.

UPDATING σ : PATH LENGTH CONTROL (CSA)

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



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

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

CSA effective and robust for $\lambda \leq N$.