

DSA4212: Natural Evolution Strategies

Bao Jiaqi, A0211257N, bao.jiaqi@u.nus.edu

April 16, 2024

1 Introduction

Natural Evolution Strategies (NES) is a family of optimization algorithms that are well-suited for black-box optimization problems. In this report, we begin by introducing and outlining optimization algorithms and their variants. We delve into their derivation and establish the connection between different algorithms. Next, we apply and compare these algorithms to three classical optimization problems. Finally, we apply the most effective method to a real-world hyperparameter tuning problem.

2 Methods

In this section, we provide an introduction to the fundamental concept of the natural gradient and Evolution Strategy. We then integrate these concepts into the notion the basic Natural Evolution Strategy. Furthermore, we introduce two extensions to enhance the basic Natural Evolution Strategy.

2.1 Natural gradient

In gradient ascent, we aim to maximize the function $J(\theta)$ at the n -th iteration using a learning rate of η . The update rule is given by:

$$\theta^{(n+1)} = \theta^{(n)} + \eta \nabla J(\theta^{(n)})$$

Specifically, $J(\theta)$ can be $\log p(\mathbf{x}|\theta)$, which equals the log-likelihood, \mathbf{x} is the data observed and θ is the model parameter to be estimated. It is important to note that gradient descent finds $\theta^{(n+1)}$ in the neighborhood of $\theta^{(n)}$ measured by Euclidean distance that minimize the objective function. However, the Euclidean distance between θ values does not serve as a suitable metric for assessing the distances between the distributions $p(\mathbf{x}|\theta)$.

In each step of the natural gradient, the steepest direction is

$$\mathbf{d} = \arg \min_{KL(p(\cdot|\theta^{(n)})||p(\cdot|\theta^{(n)}+\mathbf{d}))=\epsilon} J(\theta^{(n)} + \mathbf{d})$$

where $KL(p||q)$ denotes the Kullback–Leibler divergence from q to p . It can be derived that $KL(p(\mathbf{x}|\theta)||p(\mathbf{x}|\theta + \mathbf{d})) \approx \frac{1}{2} \mathbf{d}^T \mathbf{F}_{\theta^{(n)}} \mathbf{d}$, where $\mathbf{F}_{\theta^{(n)}} = \mathbb{E}_{\mathbf{x}}[\nabla \log(p(\mathbf{x}|\theta^{(n)})) \nabla \log(p(\mathbf{x}|\theta^{(n)}))^T]$ is the Fisher information matrix of distribution $p(\mathbf{x}|\theta^{(n)})$. Using a Lagrangian multiplier, one can derive the steepest direction: $\mathbf{d}^* = \mathbf{F}_{\theta^{(n)}}^{-1} \nabla J(\theta^{(n)})$, which is called *natural gradient*. The detailed derivations can be found in [Wen19].

One advantage of the natural gradient is its invariance under linear transformations of the coordinates, similar to Newton’s method. However, a significant drawback is the computational difficulty in calculating the inverse Fisher information matrix, particularly when the dimension of θ is large and matrix inversion introduces considerable noise.

2.2 Evolution Strategy

In many real-life optimization problems, it can be extremely challenging to calculate the gradient and Hessian of the objective function. Now, we turn our attention to the black box optimization task of maximizing $f(\mathbf{x})$, where we only have access to the function values of f and no information about $\nabla^2 f(\mathbf{x})$ or $\nabla f(\mathbf{x})$.

Evolution Strategies (ES) provide a solution to the black box optimization problem. The aim is to find a probability distribution parameterized by $\boldsymbol{\theta}$, such that the sample $\mathbf{x} \sim p(\cdot|\boldsymbol{\theta})$ corresponds to a good solution for $f(\mathbf{x})$. Our goal is to maximize the *expected fitness* under $\boldsymbol{\theta}$, given by:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p(\cdot|\boldsymbol{\theta})}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x}|\boldsymbol{\theta})d\mathbf{x}$$

We can derive $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p(\cdot|\boldsymbol{\theta})}[f(\mathbf{x})\nabla \log p(\mathbf{x}|\boldsymbol{\theta})]$ (see details in [WSG⁺14]). To optimize $J(\boldsymbol{\theta})$, we can apply a standard gradient ascent scheme, and the outlined procedure is presented in Algorithm 1.

Algorithm 1 Evolution Strategies (ES)

- 1: **Input:** objective function f , population size λ , learning rate η , initialized parameters $\boldsymbol{\theta}^{(0)}$, number of generations N .
 - 2: **for** each generation $n = 0, \dots, N - 1$ **do**
 - 3: Generate samples $\{\mathbf{x}_i\}_{i=1}^\lambda$ identically and independently from $p(\cdot|\boldsymbol{\theta}^{(n)})$
 - 4: Evaluate $f(\mathbf{x}_i)$ for each sample \mathbf{x}_i
 - 5: Compute $\nabla \log p(\mathbf{x}_i|\boldsymbol{\theta}^{(n)})$ for each sample \mathbf{x}_i
 - 6: Estimate the gradient from samples with Monte Carlo: $\hat{\nabla} J(\boldsymbol{\theta}^{(n)}) = \frac{1}{\lambda} \sum_{i=1}^\lambda \nabla \log p(\mathbf{x}_i|\boldsymbol{\theta}^{(n)})f(\mathbf{x}_i)$
 - 7: Update $\boldsymbol{\theta}^{(n+1)} = \boldsymbol{\theta}^{(n)} + \eta \hat{\nabla} J(\boldsymbol{\theta}^{(n)})$
 - 8: **end for**
-

2.3 Natural Evolution Strategy

In this subsection, we incorporate natural gradient into Evolution Strategies, referred to as Natural Evolution Strategy (NES). More specifically, we update $\boldsymbol{\theta}^{(n)}$ with the natural gradient in Step 7 of Algorithm 1.

Here, we assume that $p(\cdot|\boldsymbol{\theta})$ follows a Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix \mathbf{C} . We use the Cholesky decomposition $\mathbf{A}\mathbf{A}^T$ to represent the covariance matrix \mathbf{C} . The parameters $\boldsymbol{\theta}$ of the Gaussian distribution are expressed as a tuple $(\boldsymbol{\mu}, \mathbf{A})$. The log-likelihood is given by:

$$\log p(\mathbf{x}|\boldsymbol{\theta}) \propto -\frac{1}{2}|\mathbf{A}\mathbf{A}^T| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T(\mathbf{A}\mathbf{A}^T)^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (1)$$

We can derive the gradient of $\log p(\mathbf{x}|\boldsymbol{\theta})$ w.r.t. $\boldsymbol{\mu}$ as $\nabla_{\boldsymbol{\mu}} \log p(\mathbf{x}|\boldsymbol{\theta}) = \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})$. In addition, the gradient of $\log p(\mathbf{x}|\boldsymbol{\theta})$ w.r.t. \mathbf{A}_{ij} for each $i \leq j$ can be derived as $\nabla_{\mathbf{A}_{ij}} \log p(\mathbf{x}|\boldsymbol{\theta}) = \mathbf{R}_{ij} - \delta(i, j) \frac{1}{\mathbf{A}_{ii}}$, where $\mathbf{R} = (\mathbf{A}^{-1})^T(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1}$ and $\delta(i, j) = 1$ if $i = j$, $\delta(i, j) = 0$ otherwise; see derivation details of the Gaussian log likelihood in [YWSS09].

Moreover, in order to ensure the algorithm's invariance under rank-preserving transformations of f , the fitness $f(\mathbf{x}_i)$ in the computation of $\hat{\nabla} J(\boldsymbol{\theta}^{(n)})$ can be replaced by rank-based *fitness shaping*. Specifically, we replace the fitness values $\{f(\mathbf{x}_i)\}_{i=1}^\lambda$ with *utilities* $\{u_i\}_{i=1}^\lambda$. After sorting the samples in descending order based on $f(\mathbf{x}_i)$, u_i represents the utility for \mathbf{x}_i , and the ordering is such that $u_1 \geq \dots \geq u_\lambda$.

The combination of the formulas for the gradient of log likelihood with Algorithm 1 gives an outline of NES with fitness shaping for Gaussian search distribution, as depicted in Algorithm 2. If, in Step 8 of Algorithm 2, we omit the Fisher matrix and update the parameters solely based on the gradient, it is equivalent to the Evolution Strategy outlined in Algorithm 1 with fitness shaping for Gaussian search distribution. Additionally, for smoother estimates of gradients obtained from small samples, we employ momentum.

Algorithm 2 Natural Evolution Strategy (NES)

- 1: **Input:** objective function f , population size λ , learning rate η , momentum weight β , initialized parameters $\boldsymbol{\theta}^{(0)} = (\boldsymbol{\mu}^{(0)}, \mathbf{A}^{(0)})$, number of generations N and utilities $u_1 \geq \dots \geq u_\lambda$.
 - 2: **for** each generation $n = 0, \dots, N - 1$ **do**
 - 3: Generate $\mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I})$ independently and set $\mathbf{x}_i = \boldsymbol{\mu}^{(n)} + \mathbf{A}^{(n)} \mathbf{z}_i$ for each $i = 1, \dots, \lambda$.
 - 4: Evaluate $f(\mathbf{x}_i)$ for each sample \mathbf{x}_i and sort $\{\mathbf{x}_i, \mathbf{z}_i\}_{i=1}^\lambda$ based on $f(\mathbf{x}_i)$ in descending order.
 - 5: Compute $\nabla \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)})$ for each sample \mathbf{x}_i , where $\nabla_{\boldsymbol{\mu}^{(n)}} \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)}) = \mathbf{C}^{(n)-1}(\mathbf{x}_i - \boldsymbol{\mu}^{(n)})$ and $\nabla_{\mathbf{A}_{ij}^{(n)}} \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)}) = \mathbf{R}_{ij}^{(n)} - \delta(i, j) \frac{1}{\mathbf{A}_{ii}^{(n)}}$ for each $i \leq j$.
 - 6: Estimate the gradient from samples with Monte Carlo: $\hat{\nabla} J(\boldsymbol{\theta}^{(n)}) = \frac{1}{\lambda} \sum_{i=1}^\lambda \nabla \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)}) u_i$
 - 7: Estimate the Fisher matrix from samples with Monte Carlo: $\hat{\mathbf{F}}_{\boldsymbol{\theta}^{(n)}} = \frac{1}{\lambda} \sum_{i=1}^\lambda \nabla \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)}) \nabla \log p(\mathbf{x}_i | \boldsymbol{\theta}^{(n)})^T$.
 - 8: Update $\boldsymbol{\theta}^{(n+1)} = \boldsymbol{\theta}^{(n)} + \eta \hat{\mathbf{F}}_{\boldsymbol{\theta}^{(n)}}^{-1} \hat{\nabla} J(\boldsymbol{\theta}^{(n)}) + \beta(\boldsymbol{\theta}^{(n)} - \boldsymbol{\theta}^{(n-1)})$
 - 9: **end for**
-

Similar to the drawback of natural gradient mentioned in Section 2.1, vanilla NES requires computing the inverse of the estimated $\mathbf{F}_{\boldsymbol{\theta}^{(n)}}$. Specifically, the estimates based on generated samples may not be invertible, and a sufficiently large sample size λ is needed to obtain an accurate estimate. Additionally, the invariance under linear transformations of coordinates is partially lost when parameterizing the covariance matrix \mathbf{C} using its Cholesky decomposition \mathbf{A} .

2.4 Exponential Natural Evolution Strategies

There are several approaches to avoid the computation of the explicit $\hat{\mathbf{F}}_{\boldsymbol{\theta}^{(n)}}$ and enhance the vanilla NES algorithm. For instance, [YWSS09] suggests analytically computing the *exact* $\mathbf{F}_{\boldsymbol{\theta}^{(n)}}$ instead of estimating it from the samples. However, this approach still requires matrix inversion. In contrast, exponential Natural Evolution Strategies (xNES) [GSY⁺10] introduces an exponential parameterization of $\boldsymbol{\theta}$. In this subsection, we first explain the key innovations of xNES and subsequently provide a detailed algorithm.

Updating \mathbf{C} while ensuring it is both positive semi-definite and invariant under linear transformations of coordinates poses a challenge. One innovation of xNES is the use of the exponential map, which maps a symmetric matrix into a positive semi-definite matrix. We denote this exponential map as $\exp(\mathbf{M}) = \mathbf{U} \exp(\mathbf{D}) \mathbf{U}^T$, where $\mathbf{M} = \mathbf{U} \mathbf{D} \mathbf{U}^T$ represents the eigen decomposition of \mathbf{M} and $\exp(\mathbf{D})$ is computed by taking the exponential of the diagonal elements of \mathbf{D} . The positivity of $\exp(\mathbf{D})$ implies that $\exp(\mathbf{M})$ is positive semi-definite. As a result, \mathbf{C} can be represented as $\exp(\boldsymbol{\xi})$, where $\boldsymbol{\xi}$ is a symmetric matrix. This representation ensures that the update of $\boldsymbol{\xi}$ is invariant under linear transformations.

Another innovation of xNES is the linear transformation of the coordinate system at each iteration to a *natural* coordinate system, in which the current search distribution $p(\cdot | \boldsymbol{\theta})$ is a standard Gaussian.

Denoting $\theta = (\mu, A)$ as in NES, we introduce the following coordinate system for θ :

$$(\delta, M) \mapsto (\mu + A\delta, A \exp(M/2))$$

By doing so, the current search distribution $\mathcal{N}(\mu, AA^T)$ is encoded as a standard Gaussian with $(\delta, M) = (0, I)$. Importantly, the gradients and natural gradients of (δ, M) coincide. Consequently, the computation of the inverse Fisher Information matrix or even the Fisher Information matrix itself can be avoided.

The detailed procedure is given in Algorithm 3. For recommended parameter settings of the learning rates and utilities, please refer to [GSY⁺10].

Algorithm 3 Exponential Natural Evolution Strategies (xNES)

- 1: **Input:** objective function f , population size λ , learning rates η_μ, η_σ and η_B , initialized parameters $\theta^{(0)} = (\mu^{(0)}, A^{(0)})$, number of generations N and utilities $u_1 \geq \dots \geq u_\lambda$
 - 2: Set $\sigma^{(0)} \leftarrow \sqrt[p]{\det(A)}$ and $B^{(0)} \leftarrow A^{(0)}/\sigma^{(0)}$
 - 3: **for** each generation $n = 0, \dots, N - 1$ **do**
 - 4: Generate $z_i \sim \mathcal{N}(0, I)$ independently and set $x_i = \mu^{(n)} + A^{(n)}z_i$ for each $i = 1, \dots, \lambda$
 - 5: Evaluate $f(x_i)$ for each sample x_i and sort $\{x_i, z_i\}_{i=1}^\lambda$ based on $f(x_i)$ in descending order.
 - 6: Compute $G_\delta \leftarrow \sum_{i=1}^\lambda u_i z_i$
 - 7: Compute $G_M \leftarrow \sum_{i=1}^\lambda u_i (z_i z_i^T - I)$
 - 8: Compute gradient of step size $G_\sigma \leftarrow \text{tr}(G_M)/p$
 - 9: $G_B \leftarrow G_M - G_\sigma I$
 - 10: Update $\mu^{(n+1)} \leftarrow \mu^{(n)} + \eta_\mu \sigma^{(n)} B^{(n)} G_\delta$
 - 11: Update $\sigma^{(n+1)} \leftarrow \sigma^{(n)} \exp(\eta_\sigma/2 G_\sigma)$
 - 12: Update $B^{(n+1)} \leftarrow B^{(n)} \exp(\eta_B/2 G_B)$
 - 13: **end for**
-

2.5 Covariance Matrix Adaptation Evolution Strategy

xNES exhibits a strong connection to the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [HO01, Han16]. In this subsection, we will establish the link between xNES and CMA-ES, and subsequently provide an introduction to CMA-ES.

If $\eta_\sigma = \eta_B$, then step 11 and 12 of Algorithm 3 are reduced to $A^{(n+1)} \leftarrow A^{(n)} \exp(\eta_A/2 G_M)$ where $\eta_A = \eta_\sigma$. We have

$$A^{(n+1)} \leftarrow A^{(n)} \exp(\eta_A/2 \sum_{i=1}^\lambda u_i (z_i z_i^T - I)), C^{(n+1)} \leftarrow A^{(n)} \exp(\eta_C \sum_{i=1}^\lambda u_i (z_i z_i^T - I)) (A^{(n)})^T$$

where $\eta_A = \eta_C$. By the first order approximate of the exponential map, we could update $C^{(n+1)}$ with

$$A^{(n)} (I + \eta_C \sum_{i=1}^\lambda u_i (z_i z_i^T - I)) (A^{(n)})^T = (1 - \eta_C \sum_{i=1}^\lambda u_i) C^{(n)} + \eta_C \sum_{i=1}^\lambda u_i (x_i - \mu^{(n)}) (x_i - \mu^{(n)})^T \quad (2)$$

which coincides with the rank- μ update of covariance matrix in CMA-ES. However, the update rules in xNES follows from the principle of natural gradient ascent, but CMA-ES is based on heuristics.

We now introduce CMA-ES. In the basic setting of the evolution strategy with an isotropic Gaussian, updating the mean μ and variance σ^2 is straightforward. We update them using the sample mean and

variance of the selected samples with high fitness in the current generation. It's important to note that this basic strategy is heuristic-based and does not involve gradients or natural gradients. However, the successive generations' variances σ^2 tend to be highly correlated, restricting the exploration space. CMA-ES addresses this issue by introducing a covariance matrix, denoted as \mathbf{C} , and sampling from $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{C})$ in each generation. σ controls the scale of the distribution, often referred to as the "step size".

CMA-ES implements the rank- μ update, which accumulates information from previous generations in the update of \mathbf{C} . Since the estimation of covariance from a small sample is often unreliable due to the usual small population size λ , leveraging information from previous quantities is crucial.

While CMA-ES and xNES share similarities in the rank- μ update, a key difference is the utilization of the "evolution path" during the update of \mathbf{C} in CMA-ES. If we invert the sign of $\mathbf{x}_i - \boldsymbol{\mu}^{(n)}$ for every i , the update of \mathbf{C} in Eq. (1) will remain the same, but the update of $\boldsymbol{\mu}$ will result in the opposite direction. We refer to $\mathbf{x}_i - \boldsymbol{\mu}^{(n)}$ as a "step" since it is involved in the update of both $\boldsymbol{\mu}$ and \mathbf{C} . To preserve the information in the sign of the steps, we accumulate the update of $\boldsymbol{\mu}$ in each generation n into the evolution path, and update \mathbf{C} with the evolution path.

To facilitate the implementation of the algorithm from scratch, a fixed step size of $\sigma = 1$ is assumed for each generation, without controlling its value. The details are provided in Algorithm 1 (see 4). It is important to note that this version is a simplified adaptation of the full version presented in the tutorial [Han16] and the paper [NS24] of the CMA-ES python package, with certain parameters simplified.

Algorithm 4 Covariance Matrix Adaptation Evolution Strategy (CMA-ES)

- 1: **Input:** objective function f , population size λ , learning rates $\eta_\mu, \eta_C, \eta_{1,C}$ and $\eta_{2,C}$, evolution path damping parameter c_c , initialize $\mathbf{p}_C^{(0)} \leftarrow \mathbf{0}$, initialized parameters $(\boldsymbol{\mu}^{(0)}, \mathbf{C}^{(0)})$, number of generations N and utilities $u_1 \geq \dots \geq u_\lambda$
 - 2: **for** each generation $n = 0, \dots, N - 1$ **do**
 - 3: Generate $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}^{(n)}, \mathbf{C}^{(n)})$ independently for each $i = 1, \dots, \lambda$.
 - 4: Evaluate $f(\mathbf{x}_i)$ for each sample \mathbf{x}_i and sort $\{\mathbf{x}_i, z_i\}_{i=1}^\lambda$ based on $f(\mathbf{x}_i)$ in descending order.
 - 5: Update $\boldsymbol{\mu}^{(n+1)} \leftarrow \boldsymbol{\mu}^{(n)} + \eta_\mu \sum_{i=1}^\lambda u_i (\mathbf{x}_i - \boldsymbol{\mu}^{(n)})$
 - 6: Update evolution path $\mathbf{p}_C^{(n+1)} \leftarrow (1 - c_c) \mathbf{p}_C^{(n+1)} + \sqrt{c_c(2 - c_c)} \sum_{i=1}^\lambda u_i (\mathbf{x}_i - \boldsymbol{\mu}^{(n)})$
 - 7: Update $\mathbf{C}^{(n+1)} \leftarrow (1 - \eta_C) \mathbf{C}^{(n)}$
 - 8: Rank- μ update $\mathbf{C}^{(n+1)} \leftarrow \mathbf{C}^{(n+1)} + \eta_{1,C} \sum_{i=1}^\lambda u_i (\mathbf{x}_i - \boldsymbol{\mu}^{(n)}) (\mathbf{x}_i - \boldsymbol{\mu}^{(n)})^T$
 - 9: Update with evolution path $\mathbf{C}^{(n+1)} \leftarrow \mathbf{C}^{(n+1)} + \eta_{2,C} \mathbf{p}_C^{(n+1)} \mathbf{p}_C^{(n+1)T}$
 - 10: **end for**
-

In step 7 of Algorithm 4, if $\sum_{i=1}^\lambda u_i = 1$ and we disregard the evolutionary path (set $\eta_{2,C} = 0$), then this update rule is equivalent to Equation 2 in xNES since we simplify $\sigma^{(n)} = 1$ at each generation n . In step 5, if $\sum_{i=1}^\lambda u_i = 1$ and $\eta_\mu = 1$, then it is equivalent to updating the mean as $\boldsymbol{\mu}^{(n+1)} \leftarrow \sum_{i=1}^\lambda u_i \mathbf{x}_i$. If we select the top $\lambda/2$ samples with the highest fitness and set $u_i = \frac{1}{\lambda/2}$ if $i < \lambda/2$ and $u_i = 0$ otherwise, then $\boldsymbol{\mu}^{(n+1)}$ is the sample mean of the selected samples. If $\eta_{2,C} = 0$, $\eta_C = 0$, and $\eta_{1,C} = 1$, then $\mathbf{C}^{(n+1)}$ is the sample covariance of the selected samples. The simplified CMA-ES under the above settings explains the basic heuristics behind.

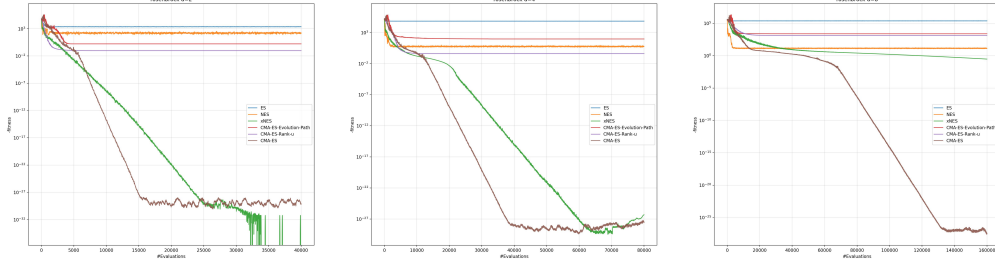


Figure 1: Rosenbrock with $d = 2, 4, 8$. We plot the evaluation of all $R(x_i)$ against the number of evaluations.

3 Applications

3.1 Experiment settings

Here are some details about the experiment settings. In ES and NES with fitness shaping for Gaussian search distribution outlined in Algorithm 2, we decay the learning rate with $\eta^{(n+1)} \leftarrow 0.9\eta^{(n)}$ in each iteration to facilitate convergence. We set the initial learning rate for μ as 0.2 and the initial learning rate for A as 0.001. We set the momentum weighted $\beta = 0.8$

Following Algorithm 4, for CMA-ES-rank- μ , we set $\eta_\mu = \eta_C = 0.2$, $\eta_{1,C} = 0.2$, and $\eta_{2,C} = 0$. For CMA-ES-Evolution-Path, we set $c_c = 0.3$, $\eta_\mu = \eta_C = 0.2$, $\eta_{1,C} = 0$, and $\eta_{2,C} = 0.3$. For CMA-ES, we set $c_c = 0.3$, $\eta_\mu = \eta_C = 0.2$, $\eta_{1,C} = 0.2$, and $\eta_{2,C} = 0.3$. For ES, NES, and CMA-ES, we set the utilities $u_i = \frac{1}{\lambda/2}$ if $i < \lambda/2$; otherwise, $u_i = 0$. For xNES, we adopt the implementation in <https://github.com/chanshing/xnes/blob/master/xnes.py> for Algorithm 3 and set $\eta_\mu = 1$, $\eta_\sigma = 0.1$, and $\eta_B = 0.01$. The population size λ is set to be proportional to the dimensionality of the function being optimized.

For the comparison of convergence, we plot the number of times we evaluate the function to maximize f against $-f(x_i)$ evaluated from generated samples.

3.2 Rosenbrock function

Rosenbrock function is a popular black-box optimization benchmark function, defined as

$$R(\mathbf{x}) = \sum_{i=1}^{d-1} [100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2]$$

We want to maximize $-R(\mathbf{x})$. We set $\lambda = 10d$ for all algorithms.

The results are compared in Figure 1. It is evident that only xNES and CMA-ES are able to successfully minimize the function, while the other methods fail to do so. Additionally, CMA-ES exhibits faster convergence compared to xNES. When d is large, only CMA-ES is effective in minimizing the Rosenbrock function. Furthermore, in scenarios with large d , NES slightly outperforms the other methods, with the exception of CMA-ES.

3.3 Ackley function

Ackley function is another popular benchmark function, defined as

$$A(\mathbf{x}) = -20 \exp \left(-0.2 \sqrt{\frac{\sum_{i=1}^d x_i^2}{d}} \right) - \exp \left(\frac{\sum_{i=1}^d \cos(2\pi x_i)}{d} \right) + 20 + \exp(1)$$

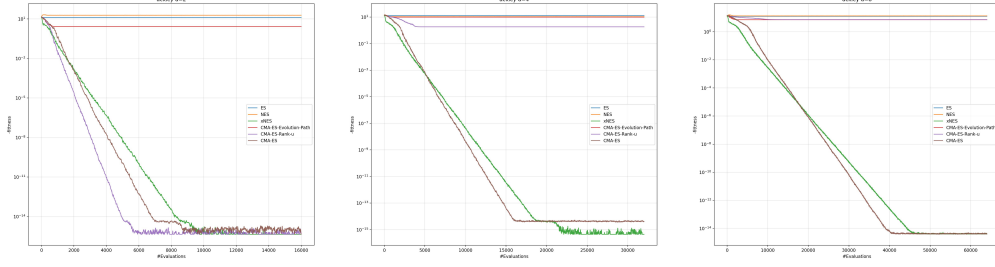


Figure 2: Ackley with $d = 2, 4, 8$. We plot the evaluation of all $A(x_i)$ against the number of evaluations.

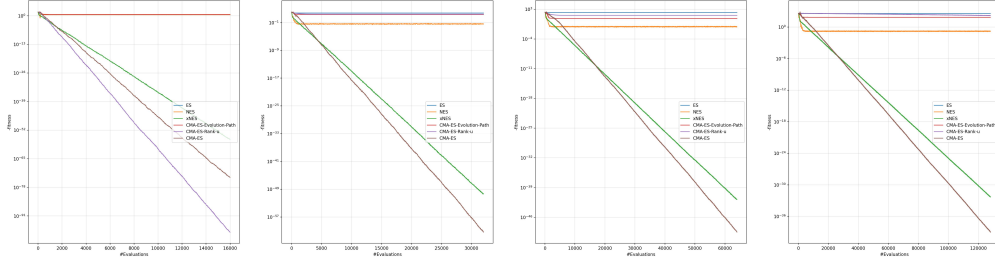


Figure 3: Sphere with $d = 2, 4, 8, 16$. We plot the evaluation of all $S(x_i)$ against the number of evaluations.

We want to maximize $-A(x)$. We set $\lambda = 4d$ for all algorithms.

The results are reported in Figure 2. It is interesting to note that when $d = 2$, CMA-ES-rank- μ exhibits the fastest convergence. However, as d increases, the integration of evolution path becomes necessary in CMA-ES-rank- μ to achieve convergence. Once again, only xNES and CMA-ES are capable of minimizing the function, with CMA-ES showing faster convergence compared to xNES.

3.4 Sphere function

Sphere function is the last benchmark function we are going to minimize, defined as $S(x) = \sum_{i=1}^d x_i^2$. We want to maximize $-S(x)$. We set $\lambda = 4d$ for all algorithms.

We report the results in Figure 3. Similarly, when $d = 2$, CMA-ES-rank- μ exhibits the fastest convergence. However, as d increases, only xNES and CMA-ES (with both rank- μ updates and evolution path) are capable of minimizing the function, with CMA-ES demonstrating faster convergence compared to xNES. Moreover, unlike the Ackley or Rosenbrock function, xNES and CMA-ES continue to minimize the sphere function. It is possible to increase the learning rate. Additionally, it is once again observed that NES yields relatively higher fitness values, particularly in relatively higher dimensions.

3.5 Hyperparameter tuning

In real-world applications, hyperparameter optimization is a challenge as it involves black box optimization. One possible solution to this is the use of evolution strategies. For our study, we utilized the CMA-ES package ¹ and Optuna ², which is a hyperparameter optimization framework that includes a built-in

¹<https://github.com/CyberAgentAILab/cmaes>

²<https://github.com/optuna/optuna>

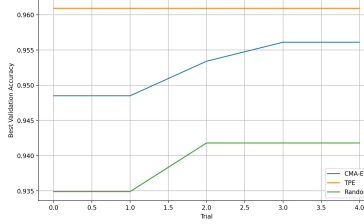


Figure 4: Optimization History. We plot best validation accuracy against number of trials.

CMA-ES sampler. In this report, we applied CMA-ES to tune a simple neural network for a multi-class classification task using the MNIST dataset ³.

The input layer has a size of 28×28 , while the output layer has a size of 10. Our goal is to optimize the size of the single hidden layer and the learning rate. We compared the performance of the CMA-ES sampler with that of the RandomSampler and Tree-structured Parzen Estimator (TPE) Sampler. TPE is a Bayesian optimization approach that samples values from the space where better outcomes are more likely.

In Figure 4, we present the comparison of the best accumulated validation accuracy against the number of trials. The results demonstrate that the TPE Sampler is the most efficient, and CMA-ES enables us to efficiently identify the best samplers compared to random search.

4 Conclusion

In this report, we have observed that xNES and CMA-ES are effective solutions for black-box optimization problems. However, it is important to note that xNES is derived from the rigorous principle of natural gradients, whereas CMA-ES is based on heuristics. Furthermore, we have found that the combination of evolution path and rank- μ updates in CMA-ES leads to improved stability and convergence. Additionally, NES outperforms ES and CMA-ES without rank- μ updates by leveraging natural gradient. Finally, we applied CMA-ES to a simple hyperparameter tuning problem and concluded that this approach is valuable in real-world scenarios.

References

- [GSY⁺10] Tobias Glasmachers, Tom Schaul, Sun Yi, Daan Wierstra, and Jürgen Schmidhuber. Exponential natural evolution strategies. In *Proceedings of the 12th annual conference on Genetic and evolutionary computation*, pages 393–400, 2010.
- [Han16] Nikolaus Hansen. The cma evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.
- [HO01] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary computation*, 9(2):159–195, 2001.
- [NS24] Masahiro Nomura and Masashi Shibata. cmaes: A simple yet practical python library for cma-es. *arXiv preprint arXiv:2402.01373*, 2024.
- [Wen19] Lilian Weng. Evolution strategies. *lilianweng.github.io*, 2019.
- [WSG⁺14] Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014.
- [YWSS09] Sun Yi, Daan Wierstra, Tom Schaul, and Jürgen Schmidhuber. Stochastic search using the natural gradient. In *Proceedings of the 26th annual international conference on machine learning*, pages 1161–1168, 2009.

³<https://pytorch.org/vision/main/generated/torchvision.datasets.MNIST.html>