Efficient Natural Evolution Strategies

Evolution Strategies and Evolutionary Programming Track

Yi Sun IDSIA Manno 6928, Switzerland yi@idsia.ch

Tom Schaul IDSIA Manno 6928, Switzerland tom@idsia.ch Daan Wierstra IDSIA Manno 6928, Switzerland daan@idsia.ch

Jürgen Schmidhuber IDSIA Manno 6928, Switzerland juergen@idsia.ch

ABSTRACT

Efficient Natural Evolution Strategies (eNES) is a novel alternative to conventional evolutionary algorithms, using the natural gradient to adapt the mutation distribution. Unlike previous methods based on natural gradients, eNES uses a fast algorithm to calculate the inverse of the exact Fisher information matrix, thus increasing both robustness and performance of its evolution gradient estimation, even in higher dimensions. Additional novel aspects of eNES include optimal fitness baselines and importance mixing (a procedure for updating the population with very few fitness evaluations). The algorithm yields competitive results on both unimodal and multimodal benchmarks.

Keywords

evolution strategies, natural gradient, optimization

1. INTRODUCTION

Evolutionary algorithms aim to optimize a 'fitness' function that is either unknown or too complex to model directly. They allow domain experts to search for good or near-optimal solutions to numerous difficult real-world problems in areas ranging from medicine and finance to control and robotics.

Typically, three objectives have to be kept in mind when developing evolutionary algorithms—we want (1) robust performance; (2) few (potentially costly) fitness evaluations; (3) scalability with problem dimensionality.

We recently introduced Natural Evolution Strategies (NES; [8]), a new class of evolutionary algorithms less ad-hoc than traditional evolutionary methods. Here we propose a novel algorithm within this framework. It retains the theoretically well-founded nature of the original NES while addressing its shortcomings w.r.t. the above objectives.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

GECCO'09, July 8–12, 2009, Montréal Québec, Canada. Copyright 2009 ACM 978-1-60558-325-9/09/07 ...\$5.00. NES algorithms maintain and iteratively update a multinormal mutation distribution. Parameters are updated by estimating a natural evolution gradient, i.e. the natural gradient on the parameters of the mutation distribution, and following it towards better expected fitness. Well-known advantages of natural gradient methods include isotropic convergence on ill-shaped fitness landscapes [2]. This avoids drawbacks of 'vanilla' (regular) gradients which are prone to slow or premature convergence [4].

Our algorithm calculates the natural evolution gradient using the *exact* Fisher information matrix (FIM) and the Monte Carlo-estimated gradient. In conjunction with the techniques of *optimal fitness baselines* and *fitness shaping* this yields robust performance (objective 1).

To reduce the number of potentially costly evaluations (objective 2), we introduce *importance mixing*, a kind of steady-state enforcer which keeps the distribution of the new population conformed to the current mutation distribution.

To keep the computational cost manageable in higher problem dimensions (objective 3), we derive a novel, efficient algorithm for computing the inverse of the exact Fisher information matrix (previous methods were either inefficient or approximate).

The resulting algorithm, Efficient Natural Evolution Strategies (eNES), is elegant, requires no additional heuristics and has few parameters that need tuning. It performs consistently well on both unimodal and multimodal benchmarks.

2. EVOLUTION GRADIENTS

First let us introduce the algorithm framework and the concept of evolution gradients. The objective is to maximize a d-dimensional unknown fitness function $f: \mathbb{R}^d \to \mathbb{R}$, while keeping the number of function evaluations – which are considered costly – as low as possible. The algorithm iteratively evaluates a population of size n individuals $\mathbf{z}_1 \dots \mathbf{z}_n$ generated from the mutation distribution $p(\mathbf{z}|\theta)$. It then uses the fitness evaluations $f(\mathbf{z}_1) \dots f(\mathbf{z}_n)$ to adjust parameters θ of the mutation distribution.

Let $J(\theta) = \mathbb{E}[f(\mathbf{z})|\theta]$ be the expected fitness under mutation distribution $p(\mathbf{z}|\theta)$, namely,

$$J(\theta) = \int f(\mathbf{z}) p(\mathbf{z}|\theta) d\mathbf{z}.$$

The core idea of our approach is to find, at each iteration, a small adjustment $\delta\theta$, such that the expected fitness

 $J\left(\theta+\delta\theta\right)$ is increased. The most straightforward approach is to set $\delta\theta\propto\nabla_{\theta}J\left(\theta\right)$, where $\nabla_{\theta}J\left(\theta\right)$ is the gradient on $J\left(\theta\right)$. Using the 'log likelihood trick', the gradient can be written as

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int f(\mathbf{z}) p(\mathbf{z}|\theta) d\mathbf{z}$$

$$= \int f(\mathbf{z}) \nabla_{\theta} p(\mathbf{z}|\theta) d\mathbf{z}$$

$$= \int f(\mathbf{z}) \frac{p(\mathbf{z}|\theta)}{p(\mathbf{z}|\theta)} \nabla_{\theta} p(\mathbf{z}|\theta) d\mathbf{z}$$

$$= \int p(\mathbf{z}|\theta) \cdot (f(\mathbf{z}) \nabla_{\theta} \ln p(\mathbf{z}|\theta)) d\mathbf{z},$$

The last term can be approximated using Monte Carlo:

$$\nabla_{\theta}^{s} J(\theta) = \frac{1}{n} \sum_{i=1}^{n} f(\mathbf{z}_{i}) \nabla_{\theta} \ln p(\mathbf{z}_{i} | \theta),$$

where $\nabla_{\theta}^{s} J(\theta)$ denotes the estimated evolution gradient.

In our algorithm, we assume that $p(\mathbf{z}|\theta)$ is a Gaussian distribution with parameters $\theta = \langle \mathbf{x}, \mathbf{A} \rangle$, where \mathbf{x} represents the mean, and \mathbf{A} represents the Cholesky decomposition of the covariance matrix \mathbf{C} , such that \mathbf{A} is upper triangular matrix and $\mathbf{C} = \mathbf{A}^{\top}\mathbf{A}$. The reason why we choose \mathbf{A} instead of \mathbf{C} as primary parameter is twofold. First, \mathbf{A} makes explicit the d(d+1)/2 independent parameters determining the covariance matrix \mathbf{C} . Second, the diagonal elements of \mathbf{A} are the square roots of the eigenvalues of \mathbf{C} , so $\mathbf{A}^{\top}\mathbf{A}$ is always positive semidefinite. In the rest of the text, we assume θ is column vector of dimension $d_s = d + d(d+1)/2$ with elements in $\langle \mathbf{x}, \mathbf{A} \rangle$ arranged as

$$\left[\left(heta^0
ight)^{ op},\left(heta^1
ight)^{ op}\ldots\left(heta^d
ight)^{ op}
ight]^{ op}.$$

Here $\theta^0 = \mathbf{x}$ and $\theta^k = [a_{k,k} \dots a_{k,d}]^\top$ for $1 \le k \le d$, where $a_{i,j}$ $(i \le j)$ denotes the (i,j)-th element of \mathbf{A} . Now we compute

$$\begin{split} \mathbf{g} \left(\mathbf{z} | \theta \right) &= & \nabla_{\theta} \ln p \left(\mathbf{z} | \theta \right) \\ &= & \nabla_{\theta} \left\{ \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{A}|^{2} \right. \\ &\left. - \frac{1}{2} \left(\mathbf{A}^{-\top} \left(\mathbf{z} - \mathbf{x} \right) \right)^{\top} \left(\mathbf{A}^{-\top} \left(\mathbf{z} - \mathbf{x} \right) \right) \right\}, \end{split}$$

where $\mathbf{g}\left(\mathbf{z}|\theta\right)$ is assumed to be a d_s -dimensional column vector. The gradient w.r.t. \mathbf{x} is simply

$$\nabla_{\mathbf{x}} \ln p(\mathbf{z}|\theta) = \mathbf{C}^{-}(\mathbf{z} - \mathbf{x}).$$

The gradient w.r.t. $a_{i,j}$ $(i \leq j)$ is given by

$$\frac{\partial}{\partial a_{i,j}} \ln p\left(\mathbf{z}|\theta\right) = r_{i,j} - \delta\left(i,j\right) a_{i,i}^{-1},$$

where $r_{i,j}$ is the (i,j)-th element of

$$\mathbf{R} = \mathbf{A}^{-\top} \left(\mathbf{z} - \mathbf{x} \right) \left(\mathbf{z} - \mathbf{x} \right)^{\top} \mathbf{C}^{-}$$

and $\delta(i, j)$ is the Kronecker Delta function.

From $\mathbf{g}(\mathbf{z}|\theta)$, the mutation gradient $\nabla_{\theta}^{s}J(\theta)$ can be computed as $\nabla_{\theta}^{s}J(\theta)=\frac{1}{n}\mathbf{G}\mathbf{f}$, where $\mathbf{G}=[\mathbf{g}(\mathbf{z}_{1}|\theta)\ldots\mathbf{g}(\mathbf{z}_{n}|\theta)]$, and $\mathbf{f}=[f(\mathbf{z}_{1})\ldots f(\mathbf{z}_{n})]^{\top}$. We update θ by $\delta\theta=\eta\nabla_{\theta}^{s}J(\theta)$, where η is an empirically tuned step size.

3. NATURAL GRADIENT

Vanilla gradient methods have been shown to converge slowly in fitness landscapes with ridges and plateaus. Natural gradients [1] constitute a principled approach for dealing with such problems. The natural gradient, unlike the vanilla gradient, has the advantage of always pointing in the direction of the steepest ascent. Furthermore, since the natural gradient is invariant w.r.t. the particular parameterization of the mutation distribution, it can cope with ill-shaped fitness landscapes and provides isotropic convergence properties, which prevents premature convergence on plateaus and avoids overaggressive steps on ridges [1].

In this paper, we consider a special case of the natural gradient $\tilde{\nabla}_{\theta} J$, defined as

$$\delta \theta^{\top} \tilde{\nabla}_{\theta} J = \max_{\delta \theta} J \left(\theta + \delta \theta \right),$$

s.t. $KL \left(\theta + \delta \theta || \theta \right) = \varepsilon,$

where ε is an arbitrarily small constant and $KL(\theta'|\theta)$ denotes the Kullback-Leibler divergence between distributions $p(\mathbf{z}|\theta')$ and $p(\mathbf{z}|\theta)$. The constraints impose a geometry on θ which differs from the plain Euclidean one. With $\varepsilon \to 0$, the natural gradient $\tilde{\nabla}_{\theta}J$ satisfies the necessary condition $\mathbf{F}\tilde{\nabla}_{\theta}J = \nabla_{\theta}J$, with \mathbf{F} being the Fisher information matrix:

$$\mathbf{F} = \mathbb{E}\left[\nabla_{\theta} \ln p \left(\mathbf{z} | \theta \right) \nabla_{\theta} \ln p \left(\mathbf{z} | \theta \right)^{\top} \right].$$

If **F** is invertible, which may not always be the case, the natural gradient can be uniquely identified by $\tilde{\nabla}_{\theta}J = \mathbf{F}^{-}\nabla_{\theta}J$, or estimated from data using $\mathbf{F}^{-}\nabla_{\theta}^{s}J$. The adjustment $\delta\theta$ can then be computed by

$$\delta\theta = \eta \mathbf{F}^- \nabla^s_{\theta} J.$$

In the following sub-sections, we show that the FIM can in fact be computed exactly, that it is invertible, and that there exists an efficient² algorithm to compute the inverse of the FIM.

3.1 Derivation of the Exact FIM

In the original NES [8], we compute the natural evolution gradient using the empirical Fisher information matrix, which is estimated from the current population. This approach has three important disadvantages. First, the empirical FIM is not guaranteed to be invertible, which could result in unstable estimations. Second, a large population size would be required to approximate the exact FIM up to a reasonable precision. Third, it is highly inefficient to invert the empirical FIM, a matrix with $O\left(d^4\right)$ elements.

We circumvent these problems by computing the exact FIM directly from mutation parameters θ , avoiding the potentially unstable and computationally costly method of estimating the empirical FIM from a population which in turn was generated from θ .

In eNES, the mutation distribution is the Gaussian defined by $\theta = \langle \mathbf{x}, \mathbf{A} \rangle$, the precise FIM **F** can be computed analytically. Namely, the (m, n)-th element in **F** is given by

$$(\mathbf{F})_{m,n} = \frac{\partial \mathbf{x}^{\top}}{\partial \theta_m} \mathbf{C}^{-} \frac{\partial \mathbf{x}}{\partial \theta_n} + \frac{1}{2} \operatorname{tr} \left(\mathbf{C}^{-} \frac{\partial \mathbf{C}}{\partial \theta_m} \mathbf{C}^{-} \frac{\partial \mathbf{C}}{\partial \theta_n} \right),$$

where θ_m , θ_n denotes the *m*-th and *n*-th element in θ . Let i_m, j_m be the a_{i_m, j_m} such that it appears at the (d+m)-th

¹For any matrix \mathbf{Q} , \mathbf{Q}^- denotes its inverse and \mathbf{Q}^\top denotes its transpose.

²Normally the FIM would involve $d_s^2 = O\left(d^4\right)$ parameters, which is intractable for most practical problems.

position in θ . First, notice that

$$\frac{\partial \mathbf{x}^{\top}}{\partial x_i} \mathbf{C}^{-} \frac{\partial \mathbf{x}}{\partial x_j} = \left(\mathbf{C}^{-} \right)_{i,j},$$

and

$$\frac{\partial \mathbf{x}^{\top}}{\partial a_{i_1,j_1}} \mathbf{C}^{-} \frac{\partial \mathbf{x}}{\partial a_{i_2,j_2}} = \frac{\partial \mathbf{x}^{\top}}{\partial x_i} \mathbf{C}^{-} \frac{\partial \mathbf{x}}{\partial a_{j,k}} = 0.$$

So the upper left corner of the FIM is C^- , and F has the following shape

$$\mathbf{F} = \left[egin{array}{cc} \mathbf{C}^- & \mathbf{0} \\ \mathbf{0} & \mathbf{F_A} \end{array}
ight].$$

The next step is to compute F_A . Note that

$$\left(\mathbf{F}_{\mathbf{A}}\right)_{m,n} = \frac{1}{2}\operatorname{tr}\left[\mathbf{C}^{-}\frac{\partial\mathbf{C}}{\partial a_{i_{m},j_{m}}}\mathbf{C}^{-}\frac{\partial\mathbf{C}}{\partial a_{i_{n},j_{n}}}\right].$$

Using the relation

$$\frac{\partial \mathbf{C}}{\partial a_{i,j}} = \frac{\partial}{\partial a_{i,j}} \mathbf{A}^{\top} \mathbf{A} = \frac{\partial \mathbf{A}^{\top}}{\partial a_{i,j}} \mathbf{A} + \mathbf{A}^{\top} \frac{\partial \mathbf{A}}{\partial a_{i,j}},$$

and the properties of the trace, we get

$$(\mathbf{F}_{\mathbf{A}})_{m,n} = \operatorname{tr} \left[\mathbf{A}^{-} \frac{\partial \mathbf{A}}{\partial a_{i_{m},j_{m}}} \mathbf{A}^{-} \frac{\partial \mathbf{A}}{\partial a_{i_{n},j_{n}}} \right]$$

$$+ \operatorname{tr} \left[\frac{\partial \mathbf{A}}{\partial a_{i_{m},j_{m}}} \mathbf{C}^{-} \frac{\partial \mathbf{A}^{\top}}{\partial a_{i_{n},j_{n}}} \right].$$

Computing the first term gives us

$$\operatorname{tr}\left[\mathbf{A}^{-}\frac{\partial\mathbf{A}}{\partial a_{i_{m},j_{m}}}\mathbf{A}^{-}\frac{\partial\mathbf{A}}{\partial a_{i_{n},j_{n}}}\right] = \left(\mathbf{A}^{-}\right)_{j_{n},i_{m}}\left(\mathbf{A}^{-}\right)_{j_{m},i_{n}}.$$

Note that since A is upper triangular, A^- is also upper triangular, so the first summand is non-zero iff

$$i_n = i_m = j_n = j_m.$$

In this case, $\left(\mathbf{A}^-\right)_{j_n,i_m}=\left(\mathbf{A}^-\right)_{j_m,i_n}=a_{j_n,i_m}^{-1},$ so

$$\operatorname{tr}\left[\mathbf{A}^{-}\frac{\partial\mathbf{A}}{\partial a_{i_{m},j_{m}}}\mathbf{A}^{-}\frac{\partial\mathbf{A}}{\partial a_{i_{n},j_{n}}}\right]=a_{i_{m},i_{n}}^{-2}\delta\left(i_{m},i_{n},j_{m},j_{n}\right).$$

Here $\delta\left(\cdot\right)$ is the generalized Kronecker Delta function, i.e. $\delta\left(i_{m},i_{n},j_{m},j_{n}\right)=1$ iff all four indices are the same. The second term is computed as

$$\operatorname{tr}\left[\frac{\partial\mathbf{A}}{\partial a_{i_{m},j_{m}}}\mathbf{C}^{-}\frac{\partial\mathbf{A}^{\top}}{\partial a_{i_{n},j_{n}}}\right]=\left(\mathbf{C}^{-}\right)_{j_{n},j_{m}}\delta\left(i_{n},i_{m}\right).$$

Therefore, we have

$$\left(\mathbf{F_{A}}\right)_{m,n} = \left(\mathbf{C}^{-}\right)_{j_{n},j_{m}} \delta\left(i_{n},i_{m}\right) + a_{i_{m},i_{n}}^{-2} \delta\left(i_{m},i_{n},j_{m},j_{n}\right).$$

It can easily be proven that $\mathbf{F}_{\mathbf{A}}$ itself is a block diagonal matrix with d blocks along the diagonal, with sizes ranging from d to 1. Therefore, the precise FIM is given by

$$\mathbf{F} = \left[egin{array}{ccc} \mathbf{F}_0 & & & & \\ & \mathbf{F}_1 & & & \\ & & \ddots & & \\ & & & \mathbf{F}_d \end{array}
ight],$$

with $\mathbf{F}_0 = \mathbf{C}^-$ and block \mathbf{F}_k $(d \ge k \ge 1)$ given by

$$\mathbf{F}_k = \left[egin{array}{cc} a_{k,k}^{-2} & \mathbf{0} \ \mathbf{0} & \mathbf{0} \end{array}
ight] + \mathbf{D}_k.$$

Here \mathbf{D}_k is the lower-right square submatrix of \mathbf{C}^- with dimension d+1-k, e.g. $\mathbf{D}_1 = \mathbf{C}^-$, and $\mathbf{D}_d = (\mathbf{C}^-)_{d,d}$.

We prove that the FIM given above is invertible if \mathbf{C} is invertible. \mathbf{F}_k $(1 \le k \le d)$ being invertible follows from the fact that the submatrix \mathbf{D}_k on the main diagonal of a positive definite matrix \mathbf{C}^- must also be positive definite, and adding $a_{k,k}^{-2} > 0$ to the diagonal would not decrease any of its eigenvalues. Also note that $\mathbf{F}_0 = \mathbf{C}^-$ is invertible, so \mathbf{F} is invertible.

It is worth pointing out that the block diagonal structure of \mathbf{F} partitions parameters θ into d+1 orthogonal groups $\theta^0 \dots \theta^k$, which suggests that we could modify each group of parameters without affecting other groups. We will need this intuition in the next section.

3.2 Iterative Computation of FIM Inverse

The exact FIM is a block diagonal matrix with d+1 blocks. Normally, inverting the FIM requires d matrix inversions. However, we can explore the structure of each sub-block in order to make the inverse of ${\bf F}$ more efficient, both in terms of time and space complexity.

First, we realize that \mathbf{F}_d is simply a number, so its inversion is given by $\mathbf{F}_d^- = \left(\left(\mathbf{C}^- \right)_{d,d} + a_{d,d}^{-2} \right)^{-1}$, and similarly $\mathbf{D}_d^- = \left(\left(\mathbf{C}^- \right)_{d,d} \right)^{-1}$. Now, letting k vary from d-1 to 1, we can compute \mathbf{F}_k^- and \mathbf{D}_k^- directly from \mathbf{D}_{k+1}^- . By block matrix inversion

$$\left[\begin{array}{cc} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{array}\right]^{-} = \left[\begin{array}{cc} \mathbf{Q}_{1}^{-} & -\mathbf{P}_{11}^{-}\mathbf{P}_{12}\mathbf{Q}_{2}^{-} \\ -\mathbf{Q}_{2}^{-}\mathbf{P}_{21}\mathbf{P}_{11}^{-} & \mathbf{Q}_{2}^{-} \end{array}\right],$$

using

$$\mathbf{Q}_1 = \mathbf{P}_{11} - \mathbf{P}_{12} \mathbf{P}_{22}^{-} \mathbf{P}_{21}, \ \mathbf{Q}_2 = \mathbf{P}_{22} - \mathbf{P}_{21} \mathbf{P}_{11}^{-} \mathbf{P}_{12},$$

and the Woodbury identity

$$\begin{aligned} \mathbf{Q}_{2}^{-} &= \left[\mathbf{P}_{22} + \mathbf{P}_{21} \left(-\mathbf{P}_{11}^{-} \right) \mathbf{P}_{21}^{\top} \right]^{-} \\ &= \mathbf{P}_{22}^{-} - \mathbf{P}_{22}^{-} \mathbf{P}_{21} \left[-\mathbf{P}_{11} + \mathbf{P}_{21}^{\top} \mathbf{P}_{22}^{-} \mathbf{P}_{21} \right]^{-} \mathbf{P}_{21}^{\top} \mathbf{P}_{22}^{-}, \end{aligned}$$

(also noting that in our case, \mathbf{P}_{11} is a number $(\mathbf{C}^-)_{k,k}$), we can state

$$\mathbf{Q}_{2}^{-} = \mathbf{P}_{22}^{-} - \frac{\left(\mathbf{P}_{22}^{-}\mathbf{P}_{21}\right)\left(\mathbf{P}_{22}^{-}\mathbf{P}_{21}\right)^{\top}}{\mathbf{P}_{21}^{+}\mathbf{P}_{22}^{-}\mathbf{P}_{21} - \mathbf{P}_{11}}.$$

This can be computed directly from \mathbf{P}_{22}^- , i.e. \mathbf{D}_{k+1}^- . Skipping the intermediate steps, we propose the following algorithm for computing \mathbf{F}_k^- and \mathbf{D}_k^- from \mathbf{D}_{k+1}^- :

$$\begin{split} \mathbf{v} &= \left(\mathbf{C}^{-}\right)_{k+1:d,k}, \, w = \left(\mathbf{C}^{-}\right)_{k,k}, \, \mathbf{u} = \mathbf{D}_{k+1}^{-}\mathbf{v}, \\ s &= \mathbf{v}^{\top}\mathbf{u}, \, q = \left(w-s\right)^{-1}, \, q_{F} = \left(w_{F}-s\right)^{-1}, \\ c &= -w^{-1}\left(1+qs\right), \, c_{F} = -w_{F}^{-1}\left(1+q_{F}s\right), \\ \mathbf{F}_{k}^{-} &= \left[\begin{array}{cc} q_{F} & c_{F}\mathbf{u}^{\top} \\ c_{F}\mathbf{u}^{\top} & \mathbf{D}_{k+1}^{-} + q_{F}\mathbf{u}\mathbf{u}^{\top} \end{array}\right], \\ \mathbf{D}_{k}^{-} &= \left[\begin{array}{cc} q & c\mathbf{u}^{\top} \\ c\mathbf{u}^{\top} & \mathbf{D}_{k+1}^{-} + q\mathbf{u}\mathbf{u}^{\top} \end{array}\right]. \end{split}$$

Here $(\mathbf{C}^-)_{k+1:d,k}$ is the sub-vector in \mathbf{C}^- at column k, and row k+1 to d. A single update from \mathbf{D}_{k+1}^- to \mathbf{F}_k^- and \mathbf{D}_k^- requires $O\left((d-k)^2\right)$ floating point multiplications. The

overall complexity of computing all sub-blocks \mathbf{F}_k^- , $1 \le k \le d$, is thus $O\left(d^3\right)$.

The algorithm is efficient both in time and storage in the sense that, on one hand, there is no explicit matrix inversion, while on the other hand, the space for storing \mathbf{D}_k (including \mathbf{F}_k , if not needed later) can be reclaimed immediately after each iteration, which means that at most $O\left(d^2\right)$ storage is required. Note also that \mathbf{F}_k^- can be used directly to compute $\delta\theta^k$, using $\delta\theta^k = \mathbf{F}_k^-\mathbf{G}^k\mathbf{f}$, where

$$\mathbf{G}^{k} = \left[\mathbf{g}^{k}\left(\mathbf{z}_{1}\right), \dots, \mathbf{g}^{k}\left(\mathbf{z}_{n}\right)\right]$$
$$= \left[\nabla_{\theta^{k}} \ln p\left(\mathbf{z}|\theta\right), \dots, \nabla_{\theta^{k}} \ln p\left(\mathbf{z}|\theta\right)\right]$$

is the submatrix of **G** w.r.t. the mutation gradient of θ^k .

To conclude, the algorithm given above efficiently computes the inverse of the exact FIM, required for computing the natural mutation gradient.

4. OPTIMAL FITNESS BASELINES

The concept of fitness baselines, first introduced in [8], constitutes an efficient variance reduction method for estimating $\delta\theta$. However, baselines as found in [5] are suboptimal w.r.t. the variance of $\delta\theta$, since this FIM may not be invertible. It is difficult to formulate the variance of $\delta\theta$ directly. However, since the exact FIM is invertible and can be computed efficiently, we can in fact compute an optimal baseline for minimizing the variance of $\delta\theta$, given by

$$\operatorname{Var}(\delta\theta) = \eta^{2} \mathbb{E}[\left(\mathbf{F}^{-} \nabla_{\theta}^{s} J - \mathbb{E}\left[\mathbf{F}^{-} \nabla_{\theta}^{s} J\right]\right)^{\top} \cdot \left(\mathbf{F}^{-} \nabla_{\theta}^{s} J - \mathbb{E}\left[\mathbf{F}^{-} \nabla_{\theta}^{s} J\right]\right)],$$

where $\triangledown^s_\theta J$ is the estimated evolution gradient, which is given by

$$\nabla_{\theta}^{s} J = \frac{1}{n} \sum_{i=1}^{n} \left[f(z_{i}) - b \right] \nabla_{\theta} \ln p(\mathbf{z}_{i} | \theta).$$

The scalar b is called the fitness baseline. Adding b does not affect the expectation of $\nabla^s_\theta J$, since

$$\mathbb{E}\left[\nabla_{\theta}^{s} J\right] = \nabla_{\theta} \int \left(f\left(\mathbf{z}\right) - b\right) p\left(\mathbf{z}|\theta\right) d\mathbf{z}$$
$$= \nabla_{\theta} \int f\left(\mathbf{z}\right) p\left(\mathbf{z}|\theta\right) d\mathbf{z}.$$

However, the variance depends on the value of b, i.e.

$$\operatorname{Var}(\delta\theta) \propto b^{2} \mathbb{E}\left[\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)^{\top}\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)\right]$$
$$-2b\mathbb{E}\left[\left(\mathbf{F}^{-}\mathbf{G}\mathbf{f}\right)^{\top}\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)\right] + \operatorname{const.}$$

Here 1 denotes a n-by-1 vector filled with 1s. The optimal value of the baseline is

$$b = \frac{\mathbb{E}\left[\left(\mathbf{F}^{-}\mathbf{G}\mathbf{f}\right)^{\top}\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)\right]}{\mathbb{E}\left[\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)^{\top}\left(\mathbf{F}^{-}\mathbf{G}\mathbf{1}\right)\right]}.$$

Assuming individuals are i.i.d., b can be approximated from data by

$$b \simeq \frac{\sum_{i=1}^{n} f(\mathbf{z}_i) \left(\mathbf{F}^{-} \mathbf{g} \left(\mathbf{z}_i \right) \right)^{\top} \left(\mathbf{F}^{-} \mathbf{g} \left(\mathbf{z}_i \right) \right)}{\sum_{i=1}^{n} \left(\mathbf{F}^{-} \mathbf{g} \left(\mathbf{z}_i \right) \right)^{\top} \left(\mathbf{F}^{-} \mathbf{g} \left(\mathbf{z}_i \right) \right)}.$$

In order to further reduce the estimation covariance, we can utilize a parameter-specific baseline for each parameter θ_j individually, which is given by

$$b_{j} = \frac{\mathbb{E}\left[\left(\mathbf{h}_{j}\mathbf{G}\mathbf{f}\right)\left(\mathbf{h}_{j}\mathbf{G}\mathbf{1}\right)\right]}{\mathbb{E}\left[\left(\mathbf{h}_{j}\mathbf{G}\mathbf{1}\right)\left(\mathbf{h}_{j}\mathbf{G}\mathbf{1}\right)\right]} \simeq \frac{\sum_{i=1}^{n} f\left(\mathbf{z}_{i}\right)\left(\mathbf{h}_{j}\mathbf{g}\left(\mathbf{z}_{i}\right)\right)^{2}}{\sum_{i=1}^{n} \left(\mathbf{h}_{j}\mathbf{g}\left(\mathbf{z}_{i}\right)\right)^{2}}.$$

Here \mathbf{h}_j is the j-th row vector of \mathbf{F}^- .

However, parameter-specific baseline values θ_j might reduce variance too much, which harms the performance of the algorithm. Additionally, adopting different baseline values for correlated parameters may affect the underlying structure of the parameter space, rendering estimations unreliable. To address both of these problems, we follow the intuition that if the (m,n)-th element in the FIM is 0, then parameters θ_m and θ_n are orthogonal and only weakly correlated. Therefore, we propose using the block fitness baseline, i.e. a single baseline b^k for each group of parameters θ^k , $0 \le k \le d$. Its formulation is given by

$$\begin{split} b^k &= \frac{\mathbb{E}\left[\left(\mathbf{F}_k^{-}\mathbf{G}^k\mathbf{f}\right)\left(\mathbf{F}_k^{-}\mathbf{G}^k\mathbf{1}\right)\right]}{\mathbb{E}\left[\left(\mathbf{F}_k^{-}\mathbf{G}^k\mathbf{1}\right)\left(\mathbf{F}_k^{-}\mathbf{G}^k\mathbf{1}\right)\right]} \\ &\simeq \frac{\sum_{i=1}^n f\left(\mathbf{z}_i\right)\left(\mathbf{F}_k^{-}\mathbf{g}^k\left(\mathbf{z}_i\right)\right)^{\top}\left(\mathbf{F}_k^{-}\mathbf{g}^k\left(\mathbf{z}_i\right)\right)}{\sum_{i=1}^n \left(\mathbf{F}_k^{-}\mathbf{g}^k\left(\mathbf{z}_i\right)\right)^{\top}\left(\mathbf{F}_k^{-}\mathbf{g}^k\left(\mathbf{z}_i\right)\right)}. \end{split}$$

Here \mathbf{F}_k^- denotes the inverse of the k-th diagonal block of \mathbf{F}^- , while \mathbf{G}^k and \mathbf{g}^k denote the submatrices corresponding to differentiation w.r.t. θ^k .

In a companion paper [7], we empirically investigated the convergence properties when using the various types of baseline. We found block fitness baselines to be very robust, whereas uniform and parameter-specific baselines sometimes led to premature convergence.

The main complexity for computing the optimal fitness baseline pertains to the necessity of storing a potentially large gradient matrix **G**, with dimension $d_s \times n \sim O(nd^2)$. The time complexity, in this case, is $O\left(nd^3\right)$ since we have to multiply each \mathbf{F}_{k}^{-} with \mathbf{G}^{k} . For large problem dimensions, the storage requirement may not be acceptable since n also scales with d. We solve this problem by introducing a time-space tradeoff which reduces the storage requirement to $O(d^2)$ while keeping the time complexity unchanged. In particular, we first compute for each k, a scalar $u_k =$ $\mathbf{a}_{k}^{-}(\mathbf{z} - \mathbf{x})$, where \mathbf{a}_{k}^{-} is the k-th row vector of \mathbf{A}^{-} . Then, for $1 \leq i \leq n$, we compute the vector $\mathbf{v} = (\mathbf{C}^{-})_{k:d}(\mathbf{z} - \mathbf{x})$, where $(\mathbf{C}^-)_{k:d}$ is the submatrix of \mathbf{C}^- by taking rows k to d. The gradient $\mathbf{g}^k(\mathbf{z}_i)$ can be computed from u_k and \mathbf{v} , and used to compute $\mathbf{F}_{k}^{-}\mathbf{g}^{k}\left(\mathbf{z}_{i}\right)$ directly. The storage for $\mathbf{g}^k(\mathbf{z}_i)$ can be immediately reclaimed. Finally, the complexity of computing $\mathbf{g}^{k}(\mathbf{z}_{i})$ for all i is O(nd(d-k)), so the total complexity of computing every element in G would still be $O(nd^3)$.

5. IMPORTANCE MIXING

At each generation, we evaluate n new individuals generated from mutation distribution $p(\mathbf{z}|\theta)$. However, since small updates ensure that the KL divergence between consecutive mutation distributions is generally small, most new individuals will fall in the high density area of the previous mutation distribution $p(\mathbf{z}|\theta')$. This leads to redundant fitness evaluations in that same area.

Our solution to this problem is a new procedure called importance mixing, which aims to reuse fitness evaluations from the previous generation, while ensuring the updated population conforms to the new mutation distribution.

Importance mixing works in two steps: In the first step, rejection sampling is performed on the previous population, such that individual z is accepted with probability

$$\min \left\{ 1, (1 - \alpha) \, \frac{p(\mathbf{z}|\theta)}{p(\mathbf{z}|\theta')} \right\}.$$

Here $\alpha \in [0,1]$ is the minimal refresh rate. Let n_a be the number of individuals accepted in the first step. In the second step, reverse rejection sampling is performed as follows: Generate individuals from $p(\mathbf{z}|\theta)$ and accept \mathbf{z} with probability

$$\max \left\{ \alpha, 1 - \frac{p(\mathbf{z}|\theta')}{p(\mathbf{z}|\theta)} \right\}$$

until $n-n_a$ new individuals are accepted. The n_a individuals from the old generation and $n-n_a$ newly accepted individuals together constitute the new population. Note that only the fitnesses of the newly accepted individuals need to be evaluated. The advantage of using importance mixing is twofold: On the one hand, we reduce the number of fitness evaluations required in each generation, on the other hand, if we fix the number of newly evaluated fitnesses, then many more fitness evaluations can potentially be used to yield more reliable and accurate gradients.

The minimal refresh rate α lower bounds the expected proportion of newly evaluated individuals $\rho = \mathbb{E}\left[\frac{n-n_{\alpha}}{\alpha}\right]$, namely $\rho \geq \alpha$, with the equality holding iff $\theta = \theta'$. In particular, if $\alpha = 1$, all individuals from the previous generation will be discarded, and if $\alpha = 0$, ρ depends only on the distance between $p(\mathbf{z}|\theta)$ and $p(\mathbf{z}|\theta')$. Normally we set α to be a small positive number, e.g. 0.01, to avoid too low an acceptance probability at the second step when $p(\mathbf{z}|\theta')/p(\mathbf{z}|\theta) \simeq 1$.

It can be proven that the updated population conforms to the mutation distribution $p(\mathbf{z}|\theta)$. In the region where $(1-\alpha) p(\mathbf{z}|\theta)/p(\mathbf{z}|\theta') \leq 1$, the probability that an individual from previous generations appears in the new population is

$$p(\mathbf{z}|\theta') \cdot (1-\alpha) p(\mathbf{z}|\theta) / p(\mathbf{z}|\theta') = (1-\alpha) p(\mathbf{z}|\theta).$$

The probability that an individual generated from the second step entering the population is $\alpha p(\mathbf{z}|\theta)$, since

$$\max \{\alpha, 1 - p(\mathbf{z}|\theta') / p(\mathbf{z}|\theta)\} = \alpha.$$

So the probability of an individual entering the population is just $p(\mathbf{z}|\theta)$ in that region. The same result holds also for the region where $(1-\alpha) p(\mathbf{z}|\theta) / p(\mathbf{z}|\theta') > 1$.

In a companion paper [7], we measured the usefulness of importance mixing, and found that it reduces the number of required fitness evaluations by a factor 5. Additionally, it reduced the algorithm's sensitivity to the population size.

The computational complexity of importance mixing can be analyzed as follows. For each generated individual \mathbf{z} , the probability $p(\mathbf{z}|\theta)$ and $p(\mathbf{z}|\theta')$ need to be evaluated, requiring $O\left(d^2\right)$ computations. The total number of individuals generated is bounded by n/α in the worst case, and is close to n on average.

6. FITNESS SHAPING

For problems with wildly fluctuating fitnesses, the gradient is disproportionately distorted by extreme fitness values,

which can lead to premature convergence or numerical instability. To overcome this problem, we use *fitness shaping*, an order-preserving nonlinear fitness transformation function [8]. The choice of (monotonically increasing) fitness shaping function is arbitrary, and should therefore be considered to be one of the tuning parameters of the algorithm. We have empirically found that ranking-based shaping functions work best for various problems. The shaping function used for all experiments in this paper was fixed to $f'(\mathbf{z}) = 2i - 1$ for i > 0.5 and $f'(\mathbf{z}) = 0$ for i < 0.5, where i denotes the relative rank of $f(\mathbf{z})$ in the population, scaled between $0 \dots 1$.

7. EFFICIENT NES

Integrating all the algorithm elements introduced above, the Efficient Natural Evolution Strategy (with block fitness baselines) can be summarized as

```
initialize \mathbf{A} \leftarrow \mathbf{I}
2
           repeat
                compute \mathbf{A}^-, and \mathbf{C}^- = \mathbf{A}^- \mathbf{A}^{-\top}
3
                 update population using importance mixing
4
5
                 evaluate f(\mathbf{z}_i) for new \mathbf{z}_i
6
                 compute rank-based fitness shaping \hat{f}
7
                 for k = d to 0
                     compute the exact FIM inverse \mathbf{F}_{k}^{-}
8
9
                     \mathbf{u} \leftarrow \mathbf{0}, \, \mathbf{v} \leftarrow \mathbf{0}, \, s_1 \leftarrow 0, \, s_2 \leftarrow 0
10
                     for i = 1 to n
                           \mathbf{q} \leftarrow \mathbf{F}_{k}^{-}\mathbf{g}^{k}\left(\mathbf{z}_{i}\right)
11
                           \mathbf{u} \leftarrow \mathbf{u} + \hat{f}(\mathbf{z}_i) \mathbf{q}

\mathbf{v} \leftarrow \mathbf{v} + \mathbf{q}
12
13
                    s_1 \leftarrow s_1 + \hat{f}(\mathbf{z}_i) \, \mathbf{q}^\top \mathbf{q}
s_2 \leftarrow s_2 + \mathbf{q}^\top \mathbf{q}
end
14
15
16
                     b^k \leftarrow s_1/s_2 \\ \delta\theta^k \leftarrow \mathbf{u} - b^k \mathbf{v}
17
18
19
20
                 \theta \leftarrow \theta + \eta \delta \theta
           {f until} stopping criteria is met
```

Note that vectors \mathbf{u} and \mathbf{v} in line 18 correspond to $\mathbf{F}_k^-\mathbf{G}^k\mathbf{f}$ and $\mathbf{F}_k^-\mathbf{G}^k\mathbf{1}$, respectively. Summing up the analysis from previous sections, the time complexity of processing a single generation is $O\left(nd^3\right)$, while the space complexity is just $O\left(d^2+nd\right)$, where $O\left(nd\right)$ comes from the need of storing the population. Assuming that n scales linearly with d, our algorithms scales linearly in space and quadratically in time w.r.t. the number of parameters, which is $O\left(d^2\right)$. This is a significant improvement over the original NES, whose complexity is $O\left(d^4\right)$ in space and $O\left(d^6\right)$ in time.

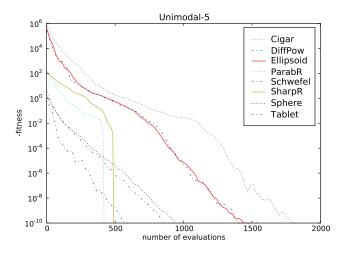
Implementations of eNES are available in both Python and Matlab³.

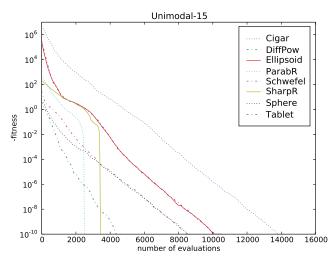
8. EXPERIMENTS

The tunable parameters of Efficient Natural Evolution Strategies are comprised of the population size n, the learning rate η , the refresh rate α and the fitness shaping function. In addition, three kinds of fitness baselines can be used.

We empirically find a good and robust choice for the learning rate η to be 1.0. On some (but not all) benchmarks the

³The Python code is part of the PyBrain machine learning library (www.pybrain.org) and the Matlab code is available at www.idsia.ch/~sun/enes.html





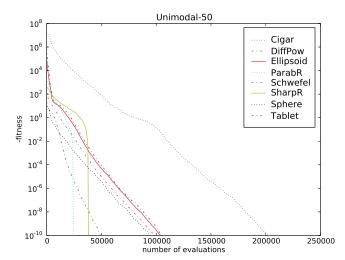


Figure 1: Results on the unimodal benchmark functions for dimension 5, 15 and 50 (from top to bottom).

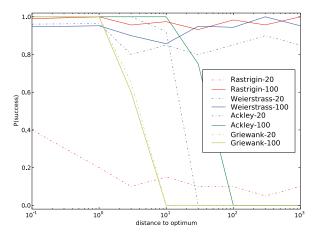


Figure 2: Success percentages varying with initial distances for the multimodal test functions using population sizes 20 and 100.

performance can be further improved by more aggressive updates. Therefore, the only parameter that needs tuning in practice is the population size, which is dependent on both the expected ruggedness of the fitness landscape and the problem dimensionality.

8.1 Benchmark Functions

We empirically validate our algorithm on 9 unimodal and 4 multimodal functions out of the set of standard benchmark functions from [6] and [3], that are typically used in the literature, for comparison purposes and for competitions. We randomly choose the inital guess at average distance 1 from the optimum. In order to prevent potentially biased results, we follow [6] and consistently transform (by a combined rotation and translation) the functions' inputs, making the variables non-separable and avoiding trivial optima (e.g. at the origin). This immediately renders many other methods virtually useless, since they cannot cope with correlated mutation directions. eNES, however, is invariant under translation and rotation. In addition, the rank-based fitness shaping makes it invariant under order-preserving transformations of the fitness function.

8.2 Performance on Benchmark Functions

We ran eNES on the set of unimodal benchmark functions with dimensions 5, 15 and 50 with population sizes 50, 250 and 1000, respectively, using $\eta=1.0$ and a target precision of 10^{-10} . Figure 1 shows the average performance over 20 runs (5 runs for dimension 50) for each benchmark function. We left out the Rosenbrock function on which eNES is one order of magnitude slower than on the other functions (e.g. 150,000 evaluations on dimension 15). Presumably this is due to the fact that the principal mutation direction is updated too slowly on complex curvatures. Note that SharpR and ParabR are unbounded functions, which explains the abrupt drop-off.

For the experiments on the multimodal benchmark functions we varied the distance of the initial guess to the optimum between 0.1 and 1000. Those runs were performed on dimension 2 with a target precision of 0.01, since here the

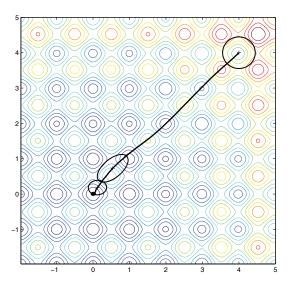


Figure 3: Evolution path and mutation distributions for a typical run on Rastrigin. Ellipsoids correspond to 0.5 standard deviations of the mutation distributions in generations 1, 20, 40.

focus was on avoiding local maxima. We compare the results for population size 20 and 100 (with $\eta=1.0$). Figure 2 shows, for all tested multimodal functions, the percentage of 100 runs where eNES found the global optimum (as opposed to it getting stuck in a local extremum) conditioned on the distance from the initial guess to the optimum.

Note that for Ackley and Griewank the success probability drops off sharply at a certain distance. For Ackley this is due to the fitness landscapes providing very little global structure to exploit, whereas for Giewank the reason is that the local optima are extremely large, which makes them virtually impossible to escape from⁴. Figure 3 shows the evolution path of a typical run on Rastrigin, and the ellipses corresponding to the mutation distribution at different generations, illustrating how eNES jumps over local optima to reach the global optimum.

For three functions we find that eNES finds the global optimum reliably, even with a population size as small as 20. For the other one, Rastrigin, the global optimum is only reliably found when using a population size of 100.

9. DISCUSSION

Unlike most evolutionary algorithms, eNES boasts a relatively clean derivation from first principles. Using a full multinormal mutation distribution and fitness shaping, the eNES algorithm is invariant under translation and rotation and under order-preserving transformations of the fitness function.

Comparing our empirical results to CMA-ES [3], considered by many to be the 'industry standard' of evolutionary computation, we find that eNES is competitive but slower, especially on higher dimensions. However, eNES is faster

on DiffPow for all dimensions. On multimodal benchmarks eNES is competitive with CMA-ES as well, as compared to the results in [8]. Our results collectively show that eNES can compete with state of the art evolutionary algorithms on standard benchmarks.

Future work will also address the problems of automatically determining good population sizes and dynamically adapting the learning rate. Moreover, we plan to investigate the possibility of combining our algorithm with other methods (e.g. Estimation of Distribution Algorithms) to accelerate the adaptation of covariance matrices, improving performance on fitness landscapes where directions of ridges and valleys change abruptly (e.g. the Rosenbrock benchmark).

10. CONCLUSION

Efficient NES is a novel alternative to conventional evolutionary algorithms, using a natural evolution gradient to adapt the mutation distribution. Unlike previous natural gradient methods, eNES quickly calculates the inverse of the exact Fisher information matrix. This increases robustness and accuracy of the evolution gradient estimation, even in higher-dimensional search spaces. Importance mixing prevents unnecessary redundancy embodied by individuals from earlier generations. eNES constitutes a competitive, theoretically well-founded and relatively simple method for artificial evolution. Good results on standard benchmarks affirm the promise of this research direction.

11. ACKNOWLEDGMENTS

This research was funded by SNF grants 200020-116674/1, 200021-111968/1 and 200021-113364/1.

12. REFERENCES

- S. Amari. Natural gradient works efficiently in learning. Neural Computation, 10(2):251–276, 1998.
- [2] S. Amari and S. C. Douglas. Why natural gradient? volume 2, pages 1213–1216 vol.2, 1998.
- [3] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.
- [4] J. Peters and S. Schaal. Natural actor-critic. Neurocomput., 71(7-9):1180-1190, 2008.
- [5] J. R. Peters. Machine learning of motor skills for robotics. PhD thesis, Los Angeles, CA, USA, 2007. Adviser-Stefan Schaal.
- [6] P. N. Suganthan, N. Hansen, J. J. Liang, K. Deb, Y. P. Chen, A. Auger, and S. Tiwari. Problem definitions and evaluation criteria for the cec 2005 special session on real-parameter optimization. Technical report, Nanyang Technological University, Singapore, 2005.
- [7] Y. Sun, D. Wierstra, T. Schaul, and J. Schmidhuber. Stochastic search using the natural gradient. In To appear in: Proceedings of the International Conference on Machine Learning (ICML-2009), 2009.
- [8] D. Wierstra, T. Schaul, J. Peters, and J. Schmidhuber. Natural evolution strategies. In Proceedings of the Congress on Evolutionary Computation (CEC08), Hongkong. IEEE Press, 2008.

 $^{^4{\}rm A}$ solution to this would be to start with a significantly larger initial C, instead of I