# Evolutionary Gradient Descent for Non-convex Optimization

Ke Xue[1], Chao Qian[1]*, Ling Xu[2], Xudong Fei[2].
{xuek, qianc}@lamda.nju.edu.cn
[1]LAMDA Group, Nanjing University, China
[2]Huawei Technologies, China

## Background and Motivation

### Non-convex optimization

- Popular in many real-world tasks
- Hard to solve
  - First order stationary point is global minima in convex optimization.
  - However, it maybe saddle point in non-convex optimization.

*How to efficiently escape saddle points and find second order stationary point is the key issue in non-convex optimization.*

### Evolutionary Algorithms

- Global convergence
- Low efficiency, especially in high dimension

### Gradient Descent

- Perform well in high dimension and large scale tasks
- Converge to local optima generally

EAs and GD each has their advantages and disadvantages.

*Can we get better algorithm for non-convex optimization by combining the merits of EAs and GD?*

## EGD Algorithm

**Algorithm 2** EGD algorithm
**Parameter**: learning rate $\eta$, population size $N$, mutation strength $\{r^{(p)}\}_{p=1}^N$, time interval $L$ for mutation, tolerance $\epsilon$, $\epsilon'$, number $T$ of iterations
**Process**:
1: Initialize the population $\{x_0^{(p)}\}_{p=1}^N$, set $i = i_{\text{mutate}} = 0$, $update^{(p)} = 1$ for $p \in [N]$;
2: **while** $i \le T$ **do**
3:  **for** $p = 1 : N$ **do**
4:   **if** $update^{(p)} = 1$ **then**
5:    **if** $\|\nabla f(x_i^{(p)})\| \le \epsilon$ and $i - i_{\text{mutate}} > L$ **then**
6:     $update^{(p)} = 0$
7:    **else**
8:     $x_{i+1}^{(p)} \leftarrow x_i^{(p)} - \eta \nabla f(x_i^{(p)})$
9:    **end if**
10:   **end if**
11:  **end for**
12:  **if** $\forall p \in [N] : update^{(p)} = 0$ **then**
13:   Apply Algorithm 3 for mutation and selection;
14:   Set $update^{(p)} = 1$ for $p \in [N]$;
15:   $i \leftarrow i + L$
16:  **end if**
17:  $i \leftarrow i + 1$
18: **end while**

*Gradient descent update or wait for mutation*

**Algorithm 3** Mutation and Selection
1: **for** $p = 1 : N$ **do**
2:  $x_{i+1}^{(p)} \leftarrow x_i^{(p)} + \xi_i^{(p)}, \xi_i^{(p)} \sim \text{Uniform}(B(0, r^{(p)}))$;
3:  $i_{\text{mutate}} \leftarrow i$;
4:  **for** $j = (i+1) : (i+L)$ **do**
5:   $x_{j+1}^{(p)} \leftarrow x_j^{(p)} - \eta \nabla f(x_j^{(p)})$
6:  **end for**
7:  **if** $f(x_{i+1+L}^{(p)}) + \epsilon' < f(x_i^{(p)})$ **then**
8:   $escape^{(p)} = 1$
9:  **else**
10:   $escape^{(p)} = 0, x_{i+1+L}^{(p)} \leftarrow x_i^{(p)}$
11:  **end if**
12: **end for**
13: $f_{\text{mean}} = \frac{1}{N} \sum_{p=1}^N f(x_{i+1+L}^{(p)})$;
14: $x_{\text{best}} = \arg\min_{x \in \{x_{i+1+L}^{(p)}\}_{p=1}^N} f(x)$;
15: **for** $p = 1 : N$ **do**
16:  **if** $escape^{(p)} = 0$ and $f(x_{i+1+L}^{(p)}) \ge f_{\text{mean}}$ **then**
17:   $x_{i+1+L}^{(p)} \leftarrow x_{\text{best}}$
18:  **end if**
19: **end for**

*Mutation and update for L iterations*

*Selection*

## Theoretical Analysis

**Theorem 2.** *Let $f$ satisfies Assumption 1. Let the parameters of EGD satisfy that $\eta = \frac{1}{\ell}$, $L = \frac{\ell}{\sqrt{\rho\epsilon}} \cdot \iota$ and $\epsilon' = \frac{1}{100\iota^3}\sqrt{\frac{\epsilon^3}{\rho}}$, where $\iota = c\log(\frac{d\ell(f(x_0)-f^*)}{\rho\epsilon\delta})$, and $c$ is an absolute constant. Then for any $\epsilon, \delta > 0$, after running*

$$\tilde{O}(\ell(f(x_0) - f^*)/\epsilon^2)$$

*iterations, EGD will find an $\epsilon$-second-order stationary point with probability at least $1 - \delta_{egd}$, where*

$$\delta_{egd} = \frac{T^*}{4L} \cdot \prod_{p=1}^N \frac{\ell\sqrt{d\epsilon}}{r^{(p)}\sqrt{\rho}} \frac{1}{\sqrt{\pi}2^\iota\iota},$$

*and $T^* = 8\max\{50\ell(f(x_0) - f^*) \cdot \iota^4, \ell(f(x_0) - f^*)\}/\epsilon^2$.*

**Remark 1.** *EGD will have more advantage over Multi-PGD,*
*(1) when the problem dimension $d$, the Lipschitz parameters $\ell$ and $\rho$ are larger, implying that the problem is more challenging;*
*(2) when the population size $N$ is larger.*

*Theorem 2 give the iterations and probability of EGD to find **$\epsilon$**-second-order stationary point.*
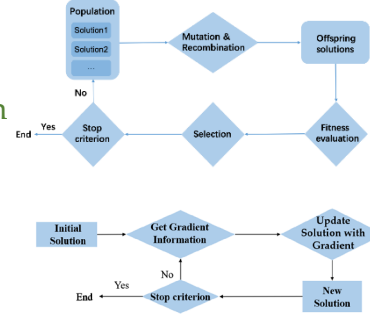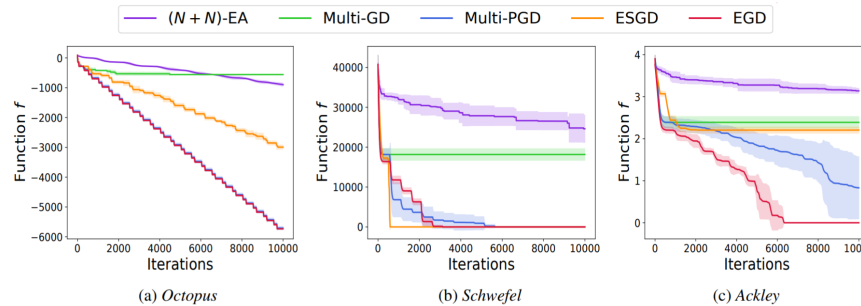
*Remark2 shows that when EGD is much better than baseline methods.*

## Experiment



(a) Octopus  (b) Schwefel  (c) Ackley

Legend: $(N + N)$-EA, Multi-GD, Multi-PGD, ESGD, EGD

| Dimension $d$ | 200 | 400 | 600 | 800 | 1000 |
|---|---|---|---|---|---|
| $f(x) < 2$ | 1.79 | 2.05 | 2.71 | 2.78 | 2.88 |
| $f(x) < 1$ | 1.70 | 2.04 | 2.24 | 2.30 | 2.34 |
| $f(x) < 0.1$ | 1.51 | 1.88 | 2.10 | 2.20 | 2.24 |

the ratio of the number of iterations of Multi-PGD and EGD until finding a solution with the value smaller than a threshold



Legend: EGD (popsize=3), EGD (popsize=5), EGD (popsize=7), EGD (popsize=10), EGD (popsize=15)



*Swimmer-v2*  *Hopper-v2*  *HalfCheetah-v2*  *Ant-v2*

Legend: Multi-GD, Multi-PGD, EGD