

Opportunistic Evolutionary Method to Minimize a Sum of Squares of Nonlinear Functions

Mikhail Zhabitsky

Joint Institute for Nuclear Research, Dubna

GRID 2016, Dubna, July 4–9

Contents

- 1 Minimization of real-parameter functions
 - Global optimization
 - Asynchronous Differential evolution
 - Crossover based on Adaptive Correlation Matrix
- 2 Non-linear Least Squares
 - Optimization problem
 - Opportunistic ADE to solve NLLS
- 3 Conclusions

Minimization of real-parameter functions

Search for a vector $\mathbf{x}^* = \{x_j\}_{j=0, \dots, D-1}$,
which minimizes an objective function $f(\mathbf{x})$:

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \Omega,$$

where $\Omega \subset \mathbb{R}^D$ is a search domain.

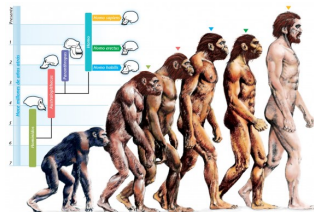
$$\mathbf{x}^* = \text{Argmin } f(\mathbf{x}); \quad f(\mathbf{x}) : \mathbb{R}^D \mapsto \mathbb{R}.$$

Additional obstacles for minimization:

- Constraints on variables $\varphi(\mathbf{x}) < 0$
- Multidimensional parameter space $D = 10 \dots 100$
- Multimodal objective functions
- Function derivatives are not available or useless
- Objective functions with “noise”
- Computationally expensive objective functions

Natural evolution

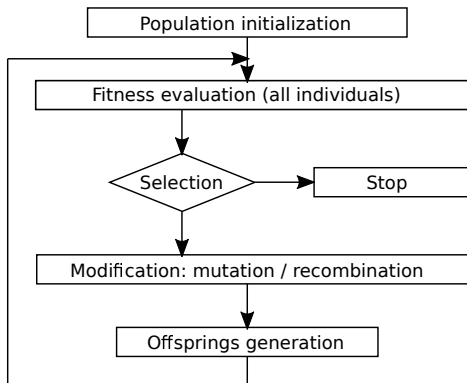
- **Natural selection** — organisms with favorable traits have a higher chance to survive and reproduce [Darwin]
- **DNA** — encoding of genetic information of the individual organism [Watson & Crick]
- **Population genetics** — modifications in offsprings as a result of mutations and recombinations and gene flow within the population [Mendel]



Result of the natural evolution — a population of individuals with increased fitness with respect to the external factors (“selection function”)

Evolutionary algorithms

Evolutionary algorithm (EA) — an algorithm, which implements operators similar to processes in Natural evolution:



Classical Differential Evolution

- Evolutionary algorithm (EA) — an algorithm, which implements operators similar to processes in Natural evolution: **selection**, **mutations** and **crossover** of genes
- Classical Differential Evolution (DE) is an Evolutionary Algorithm with specific *mutation*: $\mathbf{v} = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$
- Invented in 1995
[K. Price, R. Storn// J. Global of Optimization 11 (1997) 341]
- Use *population* of vectors (population size N_p)
- Every population member is a vector in continuous space $\Omega \subset \mathbb{R}^D$

[K. Price, R. Storn, J.A. Lampinen "Differential evolution — A Practical Approach to Global Optimization", Springer, 2005]

[S. Das, P.N. Suganthan// IEEE Trans. Evol. Comp. 15 (2011) 4]

Asynchronous Differential Evolution (ADE)

ADE is a steady-state variant of the classical DE

[E. Zhabitskaya, M. Zhabitsky // LNCS 7125 (2012) 328]

```
// population initialization  $\{\mathbf{x}_i\}_{i=0,\dots,N_p-1}$ ,  $\mathbf{x}_i = \{x_{i,j}\}_{j=0,\dots,D-1}$ 
do {
  i = choose_target_vector(); // target vector choice i
  // Mutation:
   $\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$ ; // mutant vector,  $r \neq p \neq q$  — random indices
  // Crossover (recombination):
  for ( $j = 0$ ;  $j < D$ ;  $j = j + 1$ )
    
$$u_{i,j} = \begin{cases} v_{i,j} & \text{rand}(0, 1) < C_r \text{ or } j = j_{\text{rand}} \\ x_{i,j} & \text{otherwise} \end{cases} // \text{trial vector}$$

  // Selection:
  if ( $f(\mathbf{u}_i) < f(\mathbf{x}_i)$ )
     $\mathbf{x}_i = \mathbf{u}_i$ ;
} while (termination criteria not met);
```

Asynchronous Differential Evolution: few parameters

ADE is a steady-state variant of the classical DE

[E. Zhabitskaya, M. Zhabitsky // LNCS 7125 (2012) 328]

```
// population initialization  $\{\mathbf{x}_i\}_{i=0,\dots,N_p-1}$ ,  $\mathbf{x}_i = \{x_{i,j}\}_{j=0,\dots,D-1}$ 
do {
  i = choose_target_vector(); // target vector choice i
  // Mutation:
   $\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$ ; // mutant vector,  $r \neq p \neq q$  — random indices
  // Crossover (recombination):
  for ( $j = 0$ ;  $j < D$ ;  $j = j + 1$ )
    
$$u_{i,j} = \begin{cases} v_{i,j} & \text{rand}(0, 1) < C_r \text{ or } j = j_{\text{rand}} \\ x_{i,j} & \text{otherwise} \end{cases} // \text{trial vector}$$

  // Selection:
  if ( $f(\mathbf{u}_i) < f(\mathbf{x}_i)$ )
     $\mathbf{x}_i = \mathbf{u}_i$ ;
} while (termination criteria not met);
```


Asynchronous DE: Classification

mutation:
$$\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$$

\mathbf{x}_i — target vector

\mathbf{x}_r — base vector

$(\mathbf{x}_p - \mathbf{x}_q)$ — difference vector

DE/w/x/y/z is classified according to Mutation and Crossover operators:

w corresponds to a target vector to be replaced;

x — a base vector;

y — number of difference vectors;

z — crossover type (binomial or exponential).

rand/rand/1/bin
$$\mathbf{v}_{\text{rand}} = \mathbf{x}_{\text{rand}} + F(\mathbf{x}_p - \mathbf{x}_q)$$

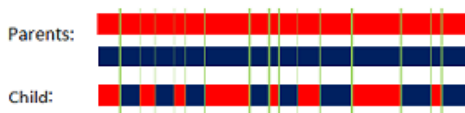
rand/best/1/bin
$$\mathbf{v}_{\text{rand}} = \mathbf{x}_{\text{best}} + F(\mathbf{x}_p - \mathbf{x}_q)$$

worst/best/bin
$$\mathbf{v}_{\text{worst}} = \mathbf{x}_{\text{best}} + F(\mathbf{x}_p - \mathbf{x}_q)$$

Uniform (Binomial) Crossover with fixed C_r

Combines a *mutant* vector \mathbf{v} with a *target* vector \mathbf{x} into a *trial* vector \mathbf{u} :

$$u_{i,j} = \begin{cases} v_{i,j} & \text{rand}(0, 1) < C_r \text{ or } j = j_{\text{rand}} \\ x_{i,j} & \text{otherwise.} \end{cases}$$



Recipes for a *crossover rate* C_r :

$C_r = 0$ is suitable for separable problems

$C_r = 1$ is for non-separable problems (also rotationally-invariant)

$C_r = 0.9$ is usually used

JADE Crossover rate C_r adaptation

JADE scheme to adapt C_r

[J. Zhang and A.C. Sanderson, IEEE TEC 13 (2009) 945]

$$C_{ri}' = N(\mu_c, \sigma_c = 0.1) \quad \text{truncated to } [0, 1]$$

Mean μ_c is updated after *successful* iterations:

$$\begin{aligned}\mu_c' &= (1 - c_c)\mu_c + c_c \langle C_{ri} \rangle \\ \mu_c^0 &= 0.5\end{aligned}$$

ADE with Adaptive Correlation Matrix

ADE-ACM: [E. Zhabitskaya, M.Zhabitsky, GECCO-2013]

DE is known to adapt to the objective function landscape:

[K. Price, R. Storn, J.A. Lampinen, Springer, 2005]

Population is a sample to test pairwise correlations:

$$q_{jk} = \frac{1}{N_p-1} \sum_{i=0}^{N_p-1} (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k),$$
$$s_{jk} = \frac{q_{jk}}{\sqrt{q_{jj}q_{kk}}},$$

Adaptive correlation matrix after *successful* iterations:

$$C' = (1 - c)C + cS.$$

[A. Auger, N. Hansen, CMA-ES, IEEE CEC (2005) 1769]

ADE with Adaptive Correlation Matrix (II)

Example: composite function

$$f = f_{\text{Rosenbrock}}(x_0, x_1) + f_{\text{Rosenbrock}}(x_2, x_3)$$

$$C = \left(\begin{array}{cc|cc} C_{00} & 0.910 & 0.027 & 0.038 \\ 0.910 & C_{11} & 0.016 & 0.015 \\ 0.027 & 0.016 & \mathbf{C_{22}} & 0.945 \\ 0.038 & 0.015 & \mathbf{0.945} & C_{33} \end{array} \right)$$

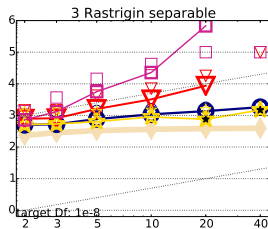
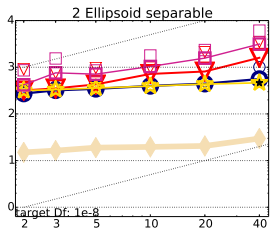
Select some coordinate m and a threshold c_{thr} :

$$m = [D\text{rand}(0, 1)], \quad c_{\text{thr}} = \text{rand}(0, 1);$$

$$\mathbb{I}_m = \{\forall j : |c_{mj}| > c_{\text{thr}}\}.$$

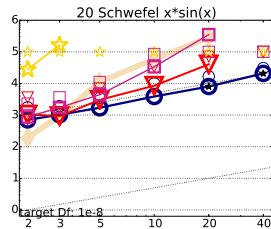
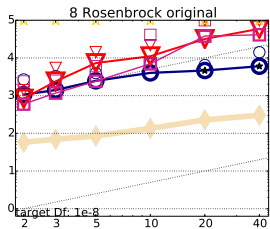
$$u_{i,j} = \begin{cases} v_{i,j} & \text{if } j \in \mathbb{I}_m \\ x_{i,j} & \text{otherwise.} \end{cases} \Rightarrow \text{Crossover within the selected subspace}$$

Crossover: numerical tests



- : ADE-ACM
- ▽: JADE
- ★: $C_r = 0.0$
- ◻: $C_r = 0.9$

rand/rand/1
 N_p adaptation



[N. Hansen et al.,
 BBOB-2012 test
 bench]

Partly-separable problems

$$D = 20 : f(\mathbf{x}) = f_{\text{Ros}}(x_0, x_1) + f_{\text{Ros}}(x_2, \dots, x_4) + \\ f_{\text{Ros}}(x_5, \dots, x_8) + f_{\text{Ros}}(x_9, \dots, x_{13}) + f_{\text{Ros}}(x_{14}, \dots, x_{19})$$

DE/rand/rand/1 strategies:

	C_r	P_{succ}	mean	median	std.dev
bin	0	0	—	—	—
bin	0.9	1	4.93e+05	4.14e+05	1.95e+05
bin	JADE	0.44	1.19e+07	—	1.11e+06
acm	—	1	6.82e+04	6.57e+04	1.74e+04

Population size N_p adaptation

	small N_p	large N_p
Probability to locate a minimum	–	+
Convergence rate	+	–

ADE with restarts [Zh&Zh// LNCS 2013]:

- Small initial population size $N_p^{\min} = 10$
- Restart with larger $N_p \leftarrow kN_p$ (after stagnation)
- Switch to N_p^{\min} if $N_p^{\max} = 20D$ was exceeded

Stagnation criteria used to solve following test problems:

- $\exists j \quad \Delta x_j < \varepsilon_x \max_i \{|x_{i,j}|\}$.
- $\Delta f < \varepsilon_f \max_{i=0, \dots, N_p-1} \{|f_i|\}$.

Scale factor F adaptation

mutation: $\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$

$F > F_{\text{thr}}$ to prevent *premature convergence*

[D. Zaharie, 2002]

[E. Zhabitskaya, LNCS 7125 (2013) 322]

JADE scheme to adapt F

[J. Zhang and A.C. Sanderson, IEEE TEC 13 (2009) 945]

$$F'_i = \text{Cauchy}(\mu_F, \sigma_F = 0.1)$$

If $F'_i < F_{\text{min}} = 0.1$, it is regenerated,

If $F'_i > F_{\text{max}} = 1.0$, then $F'_i \leftarrow F_{\text{max}}$

Location parameter μ_F is updated after *successful* iterations:

$$\mu'_F = (1 - c_F)\mu_F + c_F L_2(\{F\}) = (1 - c_F)\mu_F + c_F \frac{\sum F_i^2}{\sum F_i}$$

$$\mu_F^0 = 0.5$$

Strategies

mutation: $\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$

\mathbf{x}_i — target vector

\mathbf{x}_r — base vector

$(\mathbf{x}_p - \mathbf{x}_q)$ — difference vector

rand/rand/1

$$\mathbf{v}_i = \mathbf{x}_r + F(\mathbf{x}_p - \mathbf{x}_q)$$

rand/current-to-pbest/1

$$\mathbf{v}_i = \vec{\mathbf{x}}_i + F(\vec{\mathbf{x}}_{\text{best}}^P - \vec{\mathbf{x}}_i) + F(\vec{\mathbf{x}}_r - \vec{\mathbf{x}}_q)$$

[J. Zhang and A.C. Sanderson, JADE, IEEE TEC 13 (2009) 945]

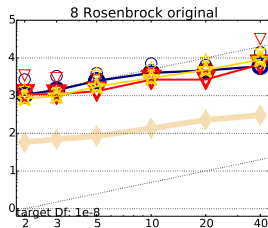
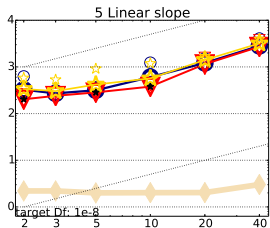
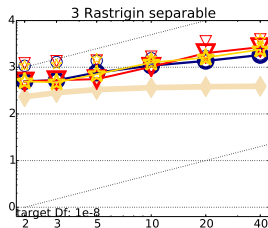
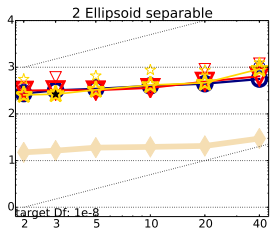
lw/current-to-pbest/1

$$\mathbf{v}_{lw} = \vec{\mathbf{x}}_{lw} + F(\vec{\mathbf{x}}_{\text{best}}^P - \vec{\mathbf{x}}_{lw}) + F(\vec{\mathbf{x}}_r - \vec{\mathbf{x}}_q)$$

$r, p, q \in P$ — random individuals from the population P

Use an archive A to produce a directional difference vector: $q \in P \cup A$

Strategies: numerical tests



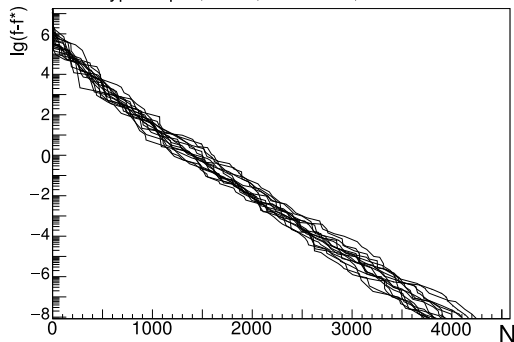
- : rand/rand/1
- ▽: rand/current-to-pbest/1
- ★: lw/current-to-pbest/1

no archive
 N_p adaptation

[N. Hansen et al.,
 BBOB-2012 test
 bench]

Convergence rate

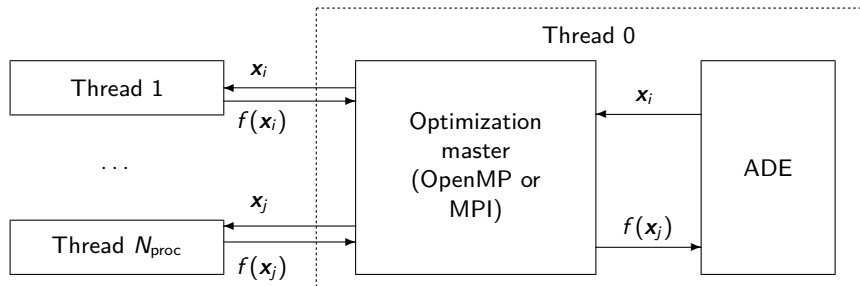
Hyperellipse, $D=10$, ADE-ACM, rand/rand/1



$$f = \sum c_j (x_j - o_j)^2$$

- log-linear convergence
- No matrix inversions
- Algorithm internal complexity $O(D^2)$

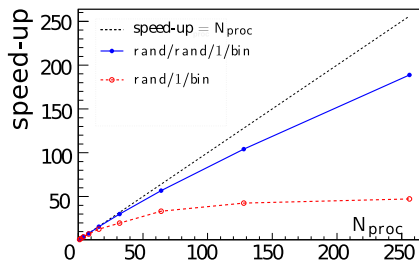
Scheme of parallelization (ADE)



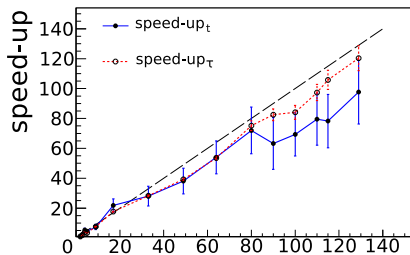
- Master/Workers model (coarse-grained parallelization)
- Complete and efficient use of all nodes
- OpenMP and MPI (MPICH2) C++ realizations

[E. Zhabitskaya, M. Zhabitsky, *Math. model.* 24 (2012) 33]

Speed-up on parallel systems



[E. Zhabitskaya, M. Zhabitsky // LNCS
 7125 (2012) 328]



[E. Zhabitskaya, et al. // Math. model.
 27 (2015) 58]

Non-linear Least Squares

Search for a vector $\mathbf{x}^* = \{x_j\}_{j=0, \dots, D-1}$,
which minimizes an objective function $f(\mathbf{x})$:

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \Omega \subset \mathbb{R}^D, \quad \text{where } f(\mathbf{x}) = \sum_{j=0}^{j < M} f_j^2(\mathbf{x}).$$
$$\mathbf{x}^* = \text{Argmin } f(\mathbf{x}); \quad f(\mathbf{x}) = \sum_{j=0}^{j < M} f_j^2(\mathbf{x}); \quad f(\mathbf{x}) : \mathbb{R}^D \mapsto \mathbb{R}^M.$$

Additional obstacles for minimization:

- **Multimodal** objective functions
- Function derivatives are not available or useless

Weighted Least Squares:
$$\chi^2 = \sum_j \frac{(y_j - t(\mathbf{x}))^2}{\sigma_j^2}$$

Note: the **Gauss–Newton** or the **Levenberg–Marquardt** algorithms for a local minimum problem

How to solve NLLS?

$$\text{Naïve example: } f(\mathbf{x}) = \sum_j f_j^2(x_j) = \sum_j c_j x_j^2.$$

Approach within Differential Evolution:

① Set $C_r = 0$ (separable problem)

② Modify a selected x_j , see improvement in $f = \sum f_j^2$

⇒ Coordinate descent: convergence rate $O(D^{-1})$

More efficient approach:

① Modify **all!** x_j , see improvements in f_j

② Combine \mathbf{x} from coordinates x_j , which provides better f_j

⇒ Convergence rate will be (almost) independent from a problem's dimension D : $O(1)$

Heuristic approach to solve NLLS

Adaptive correlation matrices for $\text{corr}(x_i, x_j)$ and $\text{corr}(x_i, f_m)$:

$$C' = (1 - c)C + cS.$$

- 1 Select a threshold $c_{\text{thr}} = \text{rand}(0, 1)$
- 2 Select some coordinate m_k and identify a set of variables
 $\mathbb{I}_k = \{\forall j : |\text{corr}(x_{m_k}, x_j)| > c_{\text{thr}}\}$
- 3 Identify a set of functions $\mathbb{F}_k = \{\forall j, l : j \in \mathbb{I}_k, |\text{corr}(x_j, f_l)| > c_{\text{thr}}\}$
- 4 Identify a set of variables $\mathbb{I}_k^{\text{ext}} = \{\forall j, l : l \in \mathbb{F}_k, |\text{corr}(x_j, f_l)| > c_{\text{thr}}\}$
- 5 Repeat steps start from 2 within $(\mathbb{I} - \mathbb{I}_k^{\text{ext}})$ until all variables are exhausted

Trial vector: $u_{i,j} = \begin{cases} v_{i,j} & \text{if } j \in \bigcup \mathbb{I}_k \\ x_{i,j} & \text{otherwise.} \end{cases} \Rightarrow$ Crossover within selected subspaces

Synthetic trial vector: $t_{i,j} = \begin{cases} u_{i,j} & \text{if } \sum_{\mathbb{F}_k} f_l(u_{i,j}) < \sum_{\mathbb{F}_k} f_l(x_{i,j}), \quad j \in \mathbb{I}_k \\ x_{i,j} & \text{otherwise.} \end{cases}$

Separable block problems: results

$$f(\mathbf{x}) = \sum_{j=0}^{D-1} f_j^2(x_j) = \sum_{j=0}^{D-1} c_j x_j^2.$$

Median number of function evaluations

D	ADE-ACM	Opportunistic ADE-ACM	Idealistic ADE-ACM
10	2486	2116	
100	76482	13365	
1000	$> 10^6$	22243	

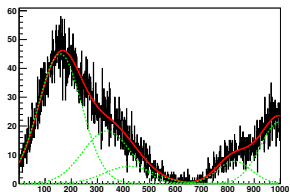
Separable block problems: results

$$f(\mathbf{x}) = \sum_{j=0}^{D-1} f_j^2(x_j) = \sum_{j=0}^{D-1} c_j x_j^2.$$

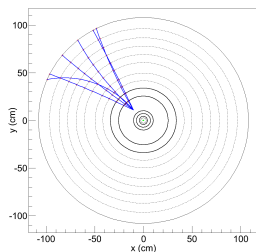
Median number of function evaluations

D	ADE-ACM	Opportunistic ADE-ACM	Idealistic ADE-ACM
10	2486	2116	596
100	76482	13365	1663
1000	$> 10^6$	22243	1893

Applications



- Curve-fitting, in particular peak-finding



- Tracking (reconstruction of particle trajectories)

Conclusions

- Opportunistic Asynchronous Differential Evolution with Adaptive Correlation Matrix is an efficient algorithm to solve partly-separable block non-linear least squares problems
- Opportunistic Asynchronous DE with ACM is (quasi) parameter-free
- Opportunistic ADE has low internal complexity $O(D^2)$
- (Almost) Linear speedup on parallel systems (up to 50 nodes)