

A0M33EOA:
EAs for Real-Parameter Optimization.
Evolution strategies. CMA-ES.

Petr Pošík

Czech Technical University in Prague
Faculty of Electrical Engineering
Department of Cybernetics

Many parts adapted or taken from
Kubalík, J. *Real-Parameter Evolutionary Algorithms*.
Lecture slides for A4M33BIA course. 2016

Introduction	2
Real EAs	3
Contents	4
Binary EAs	5
Mapping example	6
Geno-Pheno Map	7
Bit-flip mut.	8
1p xover	9
2p xover	10
Summary	11
Real EAs	12
Ops for real EAs	13
Standard ops	14
Advanced ops	16
G3	17
Summary	18
Evolution Strategies (ES)	19
Intro	20
Pipeline	21
Question	22
Gaussian Mutation	23
Adaptive Mutation	24
1/5 rule	25
Self-adaptation	27
Issues	28
CMA-ES	29
CMA-ES Demo	30
CMA-ES Code (1)	31
CMA-ES Code (2)	32
CMA-ES Code (3)	33
CMA-ES Code	34
CMA-ES Summary	35
Relations	36
Summary	37
Learning outcomes	38

EAs for real-parameter optimization

Phenotype:

- Representation that the fitness function understands and is able to evaluate.
- Vector of real numbers.

Genotype?

- Representation to which the “genetic” operators are applied.
- **Binary vector** encoding the real numbers.
 - Discretization. Finite space.
 - Discretized problem is not the same as the original one.
 - Can miss the real function optimum. Results depend on the chosen precision of discretization.
 - Requires encoding and decoding process.
- **Vector of real numbers** (genotype = phenotype).
 - *Infinite domain* (theoretically), even for space with finite bounds.
 - Opportunity to exploit *graduality* or *continuity* of the function (slight changes in variables result in slight changes of the function value).
 - No need for encoding/decoding.

Contents

Contents:

- Standard selecto-recombinative genetic algorithms with binary representation.
- Standard selecto-recombinative genetic algorithms with real representation.
- Evolution strategies.
- ES with Covariance Matrix Adaptation (CMA-ES).

Genotype-Phenotype Mapping

Assume a candidate solution to a function of 2 real numbers encoded as a single binary chromosome. Both numbers are encoded using 10 bits. The range of both numbers is $(-5.12, 5.11)$.

Which pair of real numbers corresponds to chromosome 1000000000 0100000000 ?

- A** (512, 256)
- B** (2.56, 0)
- C** ($10^9, 10^8$)
- D** (0, -2.56)

Genotype-Phenotype Mapping

Mapping binary to real vector representation (2D example):

- 2D real domain, bound constraints $[x_l, x_r] \times [y_l, y_r]$.
- Using n bits to encode each parameter.

$$x_B = [x_1 \quad x_2 \quad \dots \quad x_n \mid y_1 \quad y_2 \quad \dots \quad y_n]$$

- How to compute phenotype from known genotype?

$$x_R = \left[x_l + (x_r - x_l) \frac{\text{bin2int}(x_1, \dots, x_n)}{2^n - 1} \mid y_l + (y_r - y_l) \frac{\text{bin2int}(y_1, \dots, y_n)}{2^n - 1} \right]$$

Where in the EA should we place the mapping?

Algorithm 1: Evolutionary Algorithm

```

1 begin
2   X ← InitializePopulation()
3   f ← Evaluate(X)
4   xBSF, fBSF ← UpdateBSF(X, f)
5   while not TerminationCondition() do
6     XN ← Breed(X, f) // using certain breeding pipeline
7     fN ← Evaluate(XN)
8     xBSF, fBSF ← UpdateBSF(XN, fN)
9     X, f ← Join(X, f, XN, fN) // aka ‘replacement strategy’
10  return xBSF, fBSF

```

Algorithm 2: Evolutionary Algorithm with Genotype-Phenotype Mapping

```

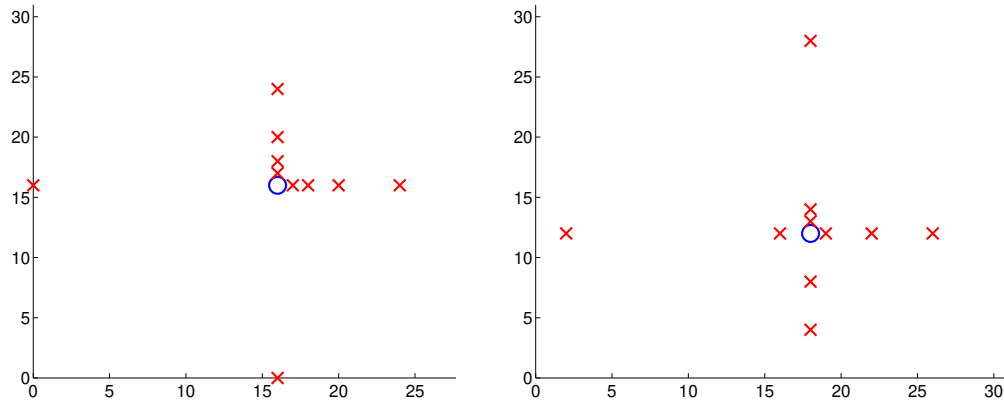
1 begin
2   X ← InitializePopulation()
3   f ← MapAndEvaluate(X)
4   xBSF, fBSF ← UpdateBSF(X, f)
5   while not TerminationCondition() do
6     XN ← Breed(X, f) // using certain breeding pipeline
7     fN ← MapAndEvaluate(XN)
8     xBSF, fBSF ← UpdateBSF(XN, fN)
9     X, f ← Join(X, f, XN, fN) // aka ‘replacement strategy’
10  return xBSF, fBSF

```

Effect of bit-flip mutation

The neighborhood of a point in the phenotype space generated by an operation applied on the genotype.

- Genotype: 10bit binary string.
- Phenotype: vector of 2 real numbers (in a discretized space).
- Operation: "bit-flip" mutation.



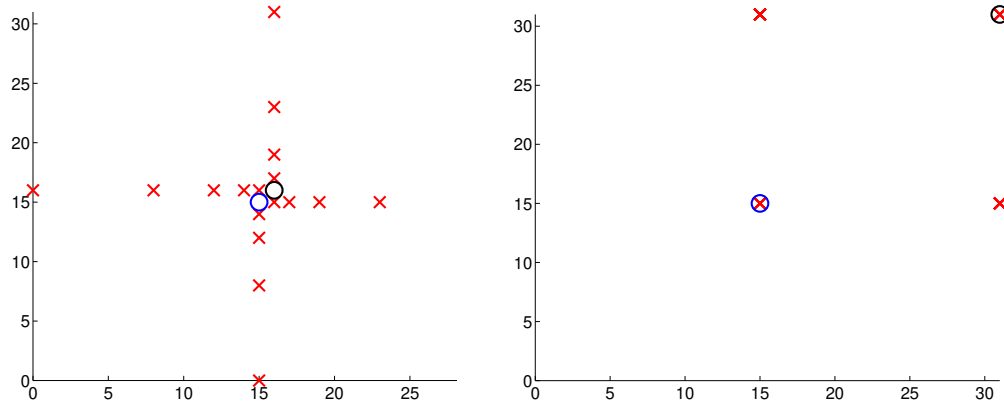
A very common situation:

- Point which is locally optimal w.r.t. the phenotype is not locally optimal w.r.t. the genotype recombination operators. (GOOD! An opportunity to escape from LO!)
- Point which is locally optimal w.r.t. the genotype recombination operators is not locally optimal w.r.t. the phenotype. (BAD: Even the best solutions found by EA do not have to correspond to the real optima we look for!)

Effect of 1-point crossover

The neighborhood of a point in the phenotype space generated by an operation applied on the genotype.

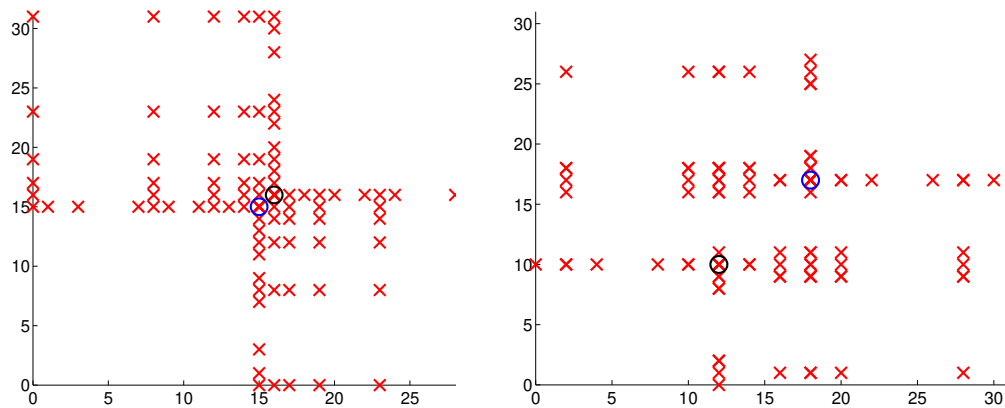
- Genotype: 10bit binary string.
- Phenotype: vector of 2 real numbers (in a discretized space).
- Operation: 1-point crossover.



Effect of 2-point crossover

The neighborhood of a point in the phenotype space generated by an operation applied on the genotype.

- Genotype: 10bit binary string.
- Phenotype: vector of 2 real numbers (in a discretized space).
- Operation: 2-point crossover.



P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 10 / 38

Summary

Binary encoding for real-parameter optimization:

- Results depend on the chosen discretization.
- The neighborhoods generated by binary crossover and mutation operators do not fit well to the “usual structures” of real-parameter functions.
- Can be useful for a rough exploration of the search space. (Then we can increase the resolution, or switch to real representation.)
- Using **Gray code** may help in certain situations, but does not solve the fundamental issues.

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 11 / 38

Recombination Operators for ESs with Real Encoding

Genotype = Phenotype = Vector of real numbers!

Standard mutation operators:

- Gaussian mutation
- Cauchy mutation

Standard recombination operators:

- Simple (1-point) Crossover: same as for binary strings
- Uniform Crossover: same as for binary strings
- Average Crossover
- Arithmetic Crossover
- Flat Crossover
- Blend Crossover BLX-(α)

Advanced recombination operators:

- Simplex Crossover (SPX)
- Unimodal Normal Distribution Crossover (UNDX)
- Parent-Centric Crossover (PCX)

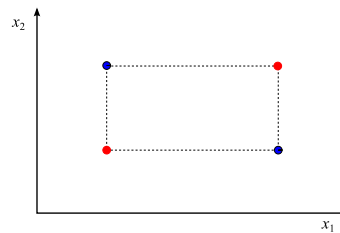
Standard Recombination Operators for Real EAs

Assume that $x^1 = (x_1^1, \dots, x_n^1)$ and $x^2 = (x_1^2, \dots, x_n^2)$ are two parents.

- **Simple (1-point) Crossover:** a position $i \in 1, 2, \dots, n - 1$ is randomly chosen, and two offspring chromosomes y^1 and y^2 are built as follows:

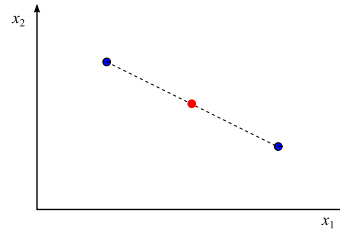
$$y^1 = (x_1^1, \dots, x_i^1, x_{i+1}^2, \dots, x_n^2)$$

$$y^2 = (x_1^2, \dots, x_i^2, x_{i+1}^1, \dots, x_n^1)$$



- **Average Crossover:** an offspring y is created as an average of the parents:

$$y = \frac{1}{2}(x^1 + x^2)$$



- **Arithmetic Crossover:** an offspring is created as a *weighted average* of the parents:

$$y = r \cdot x^1 + (1 - r) \cdot x^2,$$

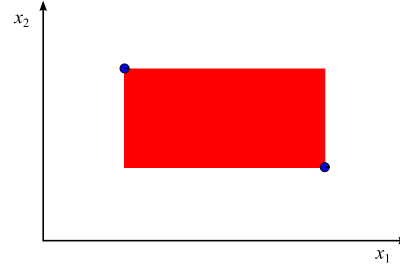
where $r \in (0, 1)$ is a constant, or varies with regard to the number of generations made, or is randomly chosen.



Standard Recombination Operators for Real EAs (cont.)

- **Flat Crossover:** an offspring $y = (y_1, \dots, y_n)$ is created such that each y_i is sampled with uniform distribution from interval

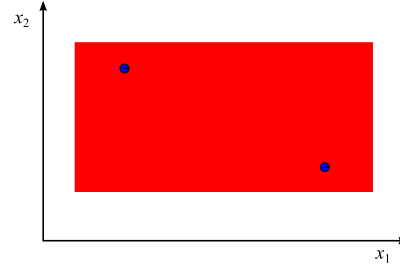
$$y_i \in [\min(x_i^1, x_i^2), \max(x_i^1, x_i^2)].$$



- **Blend Crossover:** an offspring $y = (y_1, \dots, y_n)$ is created such that each y_i is sampled with uniform distribution from interval

$$y_i \in [c_{\min} - \alpha I, c_{\max} + \alpha I],$$

where $c_{\min} = \min(p_i^1, p_i^2)$, $c_{\max} = \max(p_i^1, p_i^2)$,
 $I = c_{\max} - c_{\min}$, and $\alpha > 0$.



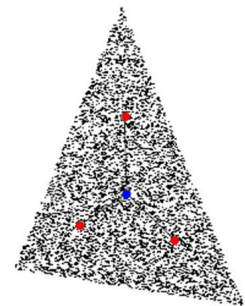
Characteristics:

- Simple, and average crossovers are deterministic; arithmetic crossover does not introduce enough diversity either.
- Simple, flat, and blend crossovers are not rotationally invariant.

Advanced Operators

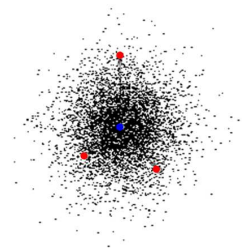
Simplex Crossover (SPX):

- Generates offspring around the mean of the μ parents
- with uniform distribution
- in a simplex which is $\sqrt{\mu + 1}$ times bigger than the parent simplex.



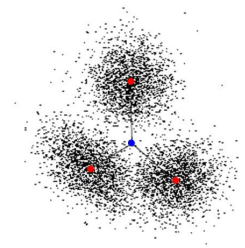
Unimodal Normal Distribution Crossover (UNDX):

- Generates offspring around the mean of the μ parents
- with multivariate normal distribution.
- Preserves the correlation among parameters well.



Parent-Centric Crossover (PCX):

- Generates offspring around one of the parents
- with multivariate normal distribution.
- The distribution shape is determined by the relative positions of the parents.
- Similar to adaptive mutation.

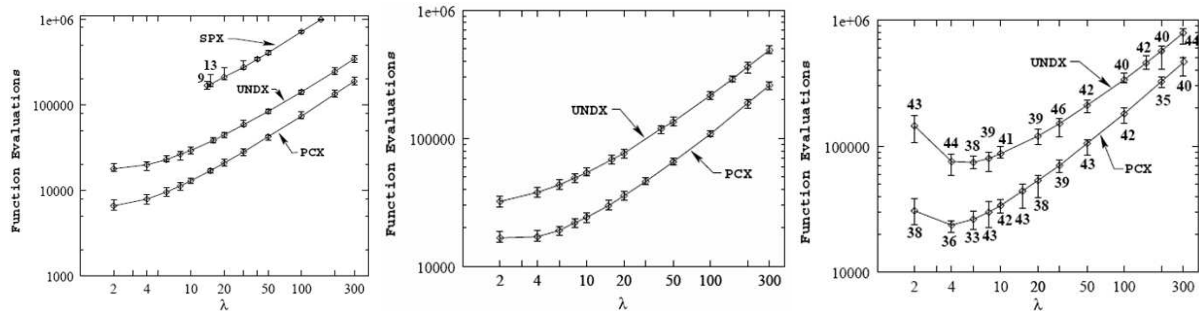


Generalized Generation Gap (G3) Algorithm

G3 [Deb05]: Elite preserving, steady-state, computationally fast. Special breeding pipeline and replacement operator.

1. From the population $P(t)$, select the best parent and $(\mu - 1)$ other parents randomly.
2. Generate λ offspring from μ parents using a recombination scheme.
3. Choose two parents at random from μ parents.
4. Form a combined subpopulation of chosen two parents and λ offspring, choose the best two solutions and replace the chosen two parents with these solutions.

Comparisons of UNDX, SPX and PCX with the G3 model on Ellipsoidal, Schwefel's, and Generalized Rosenbrock's functions for $D = 20$.



[Deb05] K. Deb. A population-based algorithm-generator for real-parameter optimization. *Soft Computing*, 9(4):236–253, April 2005.

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 17 / 38

Summary

Selecto-recombinative standard EAs with real encoding

- often use the same algorithm and breeding pipeline as binary EAs,
- although a specialized pipeline can be designed (e.g., G3).
- They use different mutation and crossover operators.

Operators for real encoding:

- Much wider range of possibilities than in binary space.
- Generally, there is no single best operator for all problems.
- Operators resulting in normal distribution of offspring usually work better for practical problems.

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 18 / 38

Evolution Strategies: Introduction

“The European branch of Evolutionary Computation.”

- Originated in Germany in 1960’s (Ingo Rechenberg and Hans-Paul Schwefel).
- ES use the natural representation of vectors in R^D as “chromosomes”.
- ES originally relied on *mutation and selection* only; recombination was added later.
- Mutation is performed by adding a random vector distributed according to multivariate Gaussian with covariance matrix $\sigma\mathbf{I}$, $\text{diag}(\sigma_1, \dots, \sigma_D)$, or general C .
- Special feature: *built-in* adaptation of mutation parameters!

Notation: $(\mu \dagger \lambda)$ -ES

- μ is the *population size* (and number of parents),
- λ is the *number of offspring* created each generation,
- $+ \text{ or } ,$ denote the *replacement strategy*:
 - $,$ is *generational* strategy: old population is discarded, new population of μ parents is chosen from the λ generated offspring.
 - $+$ is *steady-state* strategy: old population is joined with the new offspring, new population of μ parents is chosen from the joined $\mu + \lambda$ individuals.

Notation: $(\mu/\rho \dagger \lambda)$ -ES

- *Recombination* (usually deterministic), choose ρ individuals out of μ parents, $\mu \geq \rho$.
- Sometimes, subscript to ρ is used to denote the type of recombination, e.g., ρ_I for intermediate recombination (average), or ρ_W for weighted recombination (weighted average). Other recomb. ops from Real EAs can be used in principle.

Evolution Strategy Algorithm

ES use ordinary EA template (see lecture 1), with only slightly changed pipeline:

Algorithm 3: ES Breeding Pipeline

Input: Population X of μ individuals, with their fitness in f .

Number of parents ρ . Number of offspring λ .

Output: Population X_N of λ offspring.

```

1 begin
2    $X_N \leftarrow \emptyset$ 
3   for  $i \leftarrow 1, \dots, \lambda$  do
4      $X_S \leftarrow \text{SelectParents}(X, f)$  //  $\rho$  parents
5      $x_R \leftarrow \text{Recombine}(X_S)$  // usually only single offspring
6      $x_N \leftarrow \text{Mutate}(x_R)$ 
7      $X_N \leftarrow X_N \cup \{x_N\}$ 
8   return  $X_N$ 

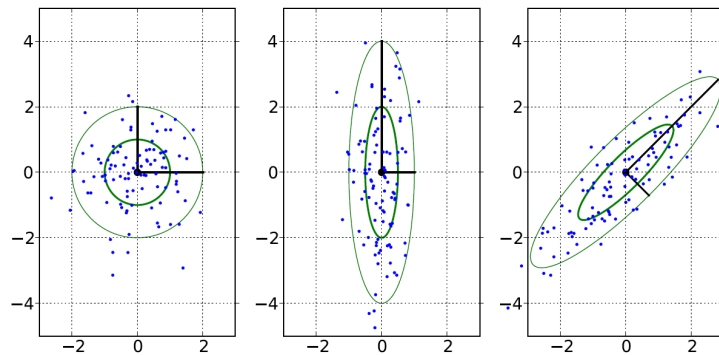
```

- The `join()` operation then forms new population for the next generation by choosing the best μ individuals either from X_N (comma strategy) or from $X \cup X_N$ (plus strategy).
- Very often $\rho = \mu$, resulting in $(\mu/\mu \dagger \lambda)$ – ES. All offspring are then centered around a single vector x_R . Lines 4 and 5 can thus be removed from the `for`-loop and placed before it.

Question: Gaussian mutation

Which of the following figures depict Gaussian distributions with covariance matrices given by

- $C = \sigma^2 I$ and
- $C = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$?



- A left and middle
- B left and right
- C middle and right
- D right and left

Gaussian Mutation

Gaussian mutation: the mutated offspring y are distributed around the original individual x as

$$y \sim N(x, C) \sim x + N(0, C) \sim x + C^{\frac{1}{2}} N(0, I),$$

where $N(\mu, C)$ is a **multivariate Gaussian distribution** with probability density function in R^D

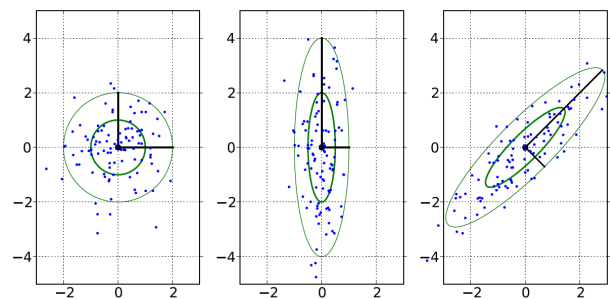
$$f_D(x|\mu, C) = \frac{1}{\sqrt{(2\pi)^D \det(C)}} \exp\left(-\frac{1}{2}(x - \mu)^T C^{-1}(x - \mu)\right)$$

Parameters:

- μ : location of the distribution. When used for mutation, $\mu = 0$ to prevent bias.
- C : Covariance matrix; determines the shape of the distribution:
 - **Isotropic:** $C = \sigma^2 I$
 - **Axis-parallel:** $C = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$
 - **General:** C positive definite

How many degrees of freedom (free parameters) do these have?

How to set up the parameters of covariance matrix?

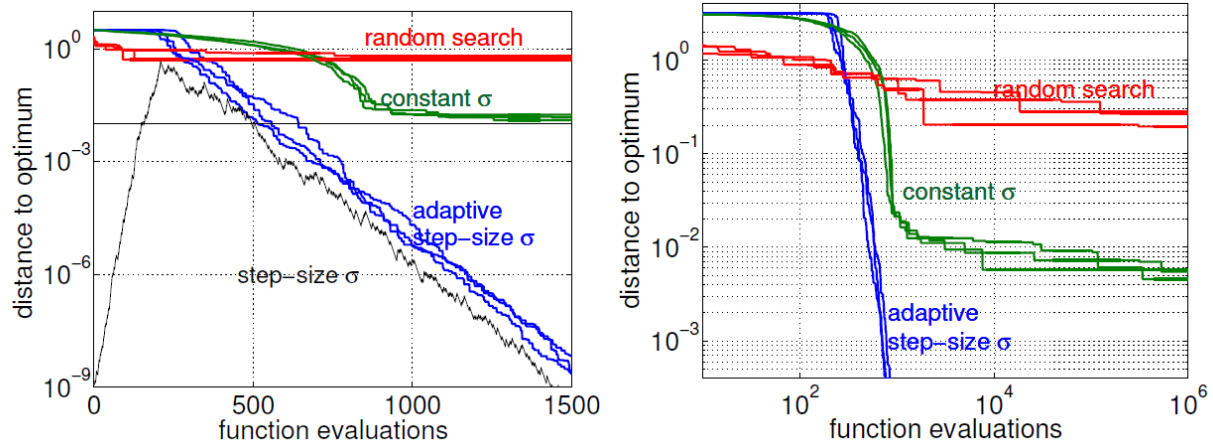


Adaptation of Mutation Parameters

Adaptation of mutation parameters is **key to ES design!**

Example: (1 + 1)-ES (hill-climber) with isotropic mutation on Sphere function: $f = \sum_i x_i^2$

- Random search vs
- (1 + 1)-ES with constant $\sigma = 10^{-2}$ vs
- (1 + 1)-ES with σ adapted using $\frac{1}{5}$ -rule with $\sigma_0 = 10^{-9}$



- Random search: inefficient.
- Constant σ : initially too small value, appropriate value between 600 and 800 evals, too large value at the end.
- Adaptive σ : near-optimal value during (almost) the whole run!

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 24 / 38

1/5 Success Rule

Progress rate φ : a ratio of the distance covered towards the optimum and the number of evaluations required to reach this distance.

Rechenberg analyzed the behavior of (1+1)-ES on 2 simple functions:

- Corridor function: $f_1(x) = x_1$ if $|x_i| < 1$ for $i \in (2, \dots, D)$, otherwise $f_1(x) = \infty$
- Sphere function: $f_2(x) = \sum_i x_i^2$

Findings:

- In both cases, the optimal step size σ^{opt} is inversely proportional to the dimension of the space D (number of variables).
- The maximum progress rate φ^{max} is also inversely proportional to D .
- For the optimal step sizes, the following probabilities of a successful mutation were obtained:
 - $p_{S,1}^{opt} = 1/(2e) \approx 0.184$
 - $p_{S,2}^{opt} \approx 0.270$

1/5 success rule: To obtain nearly optimal (local) performance of the (1+1)-ES in real-valued search spaces, tune the mutation step in such a way that the (measured) success rate is about 1/5.

- If it is greater than 1/5, increase the mutation step σ ; if it is less, decrease σ .

In practice, the 1/5 success rule has been mostly superseded by more sophisticated methods. However, its conceptual insight remain remarkably valuable.

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 25 / 38

(1+1)-ES with 1/5 rule

Algorithm 4: (1+1)-ES with 1/5 rule

```
Input:  $D \in \mathbb{N}^+, d \approx \sqrt{D+1}$ 
1 begin
2    $x \leftarrow \text{Initialize}()$ 
3   while not TerminationCondition() do
4      $x_N \leftarrow x + \sigma \mathcal{N}(0, I)$  // mutation/perturbation
5      $b \leftarrow \text{BetterThan}(x_N, x)$  // Mutation successful?
6      $\sigma \leftarrow \sigma \left( \exp \left( \mathbb{1}(b) - \frac{1}{5} \right) \right)^{\frac{1}{d}}$  // 1/5 rule
7     if  $b$  then
8        $x \leftarrow x_N$ 
```

- $\mathbb{1}(b)$ is an indicator function:

$$\mathbb{1}(b) = \begin{cases} 1 & \text{iff } b \text{ is true,} \\ 0 & \text{iff } b \text{ is false.} \end{cases}$$

- In $(1 + \lambda)$ -ES, replace b with the proportion r of improving mutations in generation.
- Other implementations are possible.

Self-adaptation

Self-adaptation:

- **Strategy parameters are part of the chromosome!** $x = (x_1, \dots, x_D, \sigma_1, \dots, \sigma_D)$
- Parameters undergo evolution together with the decision variables.
- Each individual holds information how it shall be mutated.

Example: assuming axis-parallel normal distribution is used,

- mutation of $x = (x_1, \dots, x_D, \sigma_1, \dots, \sigma_D)$ creates an offspring individual

$$x' = (x'_1, \dots, x'_D, \sigma'_1, \dots, \sigma'_D)$$

by mutating each part in a different way:

$$\sigma'_i \leftarrow \sigma_i \cdot \exp(\tau \cdot \mathcal{N}(0, 1)) \quad x'_i \leftarrow x_i + \sigma'_i \cdot \mathcal{N}(0, 1)$$

- Intuition: a “bad” σ' probably generates bad x' and is eliminated by selection.

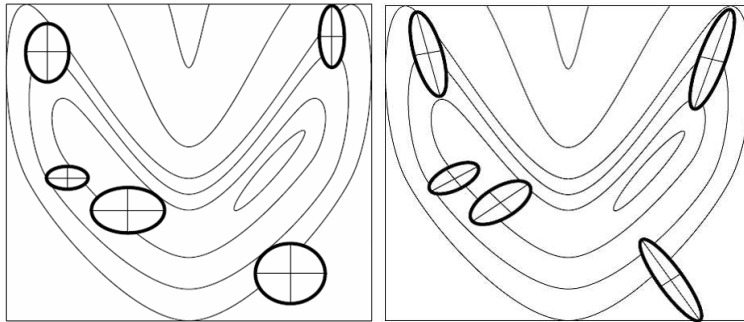
Remarks:

- An algorithm can adapt a global step size σ and coordinate-wise step sizes separately, such that the resulting coordinate-wise st. dev. is given as $\sigma \cdot \sigma_i$.
- The global step size may be adapted e.g. by the 1/5-rule.

Generalizations and issues

Generalizing from

- axis-parallel mutation distributions with D strategy parameters to
- general normal mutation distributions with full cov. matrix requires adaptation of $\frac{1}{2}D(D+1)$ strategy parameters!



Issues with self-adaptation: **selection noise** (the more parameters, the worse)!

- The intuition from the previous slide does not work much!
- A good offspring may be generated with poor strategy parameter settings (**poor setting survives**), or a bad offspring may be generated with good parameter settings (**good setting is eliminated**).

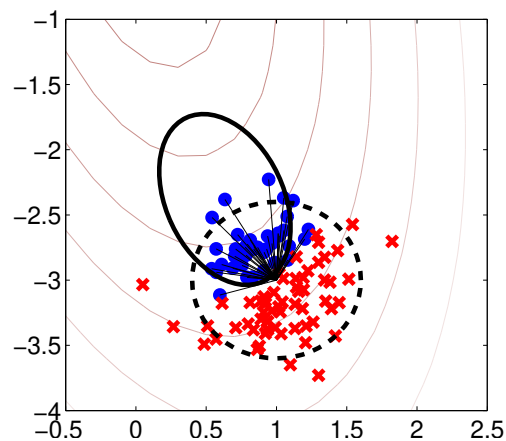
Solutions: derandomization via

- reducing the number of mutation distribution: $(1, \lambda)$ -ES, $(\mu / \mu, \lambda)$ -ES, and
- accumulating info in time (evolution paths).

CMA-ES

Evolutionary strategy with covariance matrix adaptation [HO01]:

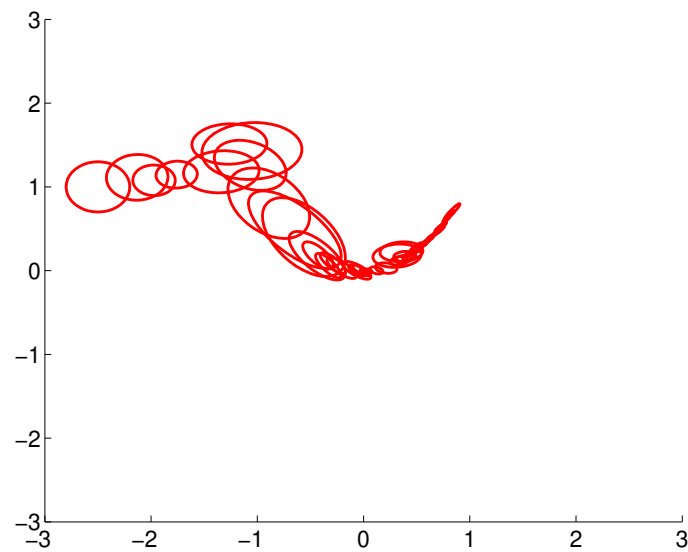
- Currently, *de facto* standard in real-parameter optimization.
- $(\mu / \mu_W, \lambda)$ -ES: recombinative, mean-centric
- Offspring is created by sampling from a single normal distribution.
- Successful mutation steps are used to adapt the mean x and the covariance matrix C of the distribution.
- Accumulates the successful steps over many generations.



[HO01] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.

CMA-ES Demo

CMA-ES on the Rosenbrock function:



P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 30 / 38

CMA-ES Code (1)

CMA-ES is a **complex, but carefully designed** and tuned algorithm!

Really? It does not seem so from the pseudocode below...

Algorithm 5: CMA-ES

```
1 begin
2   Initialize:  $x \in \mathbb{R}^D, \sigma \in \mathbb{R}_+^D, C = I$ .
3   while not TerminationCondition() do
4      $\mathcal{M} \leftarrow \text{SampleDistribution}(\lambda, \mathcal{N}(x, \sigma^2 C))$ 
5      $\mathcal{P} \leftarrow \text{SelectBest}(\mu, \mathcal{M})$ 
6      $(x, \sigma, C) \leftarrow \text{UpdateModel}(x, \sigma, C, \mathcal{P})$ 
7   return  $x$ 
```

Hm, ok, how is the Normal distribution actually sampled?

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 31 / 38

CMA-ES Code (2)

CMA-ES with the distribution sampling step expanded:

Algorithm 6: CMA-ES

```
1 begin
2   Initialize:  $x \in \mathbb{R}^D, \sigma \in \mathbb{R}_+^D, C = I$ .
3   while not TerminationCondition() do
4     for  $k \in 1, \dots, \lambda$  do
5        $z_k \leftarrow \mathcal{N}(0, I)$ 
6        $x_k \leftarrow x + \sigma C^{\frac{1}{2}} \times z_k$ 
7        $\mathcal{P} \leftarrow \text{SelectBest}(\mu, \{z_k, f(x_k) \mid 1 \leq k \leq \lambda\})$ 
8        $(x, \sigma, C) \leftarrow \text{UpdateModel}(x, \sigma, C, \mathcal{P})$ 
9   return  $x$ 
```

Remarks:

- All individuals exist in 2 “versions”: z_k distributed as $\mathcal{N}(0, I)$, and x_k distributed as $\mathcal{N}(x, \sigma^2 C)$.
- x_k are used just as an intermediate step for evaluation!
- z_k are used for model update via the population of selected parents \mathcal{P} .

OK, that’s not that complex. What about the model update?

CMA-ES Code (3)

CMA-ES with the model update step expanded:

Algorithm 7: CMA-ES

```
1 begin
2   Initialize:  $x \in \mathbb{R}^D, \sigma \in \mathbb{R}_+^D, C = I, s_\sigma = 0, s_c = 0$ .
3   while not TerminationCondition() do
4     for  $k \in 1, \dots, \lambda$  do
5        $z_k \leftarrow \mathcal{N}(0, I)$ 
6        $x_k \leftarrow x + \sigma C^{\frac{1}{2}} \times z_k$ 
7        $\mathcal{P} \leftarrow \text{SelectBest}(\mu, \{z_k, f(x_k) \mid 1 \leq k \leq \lambda\})$ 
8        $s_\sigma \leftarrow (1 - c_\sigma) s_\sigma + \sqrt{c_\sigma(2 - c_\sigma)} \sqrt{\mu_w} \sum_{z_k \in \mathcal{P}} w_k z_k$  // search path for  $\sigma$ 
9        $s_c \leftarrow (1 - c_c) s_c + h_\sigma \sqrt{c_c(2 - c_c)} \sqrt{\mu_w} \sum_{z_k \in \mathcal{P}} w_k C^{\frac{1}{2}} z_k$  // search path for  $C$ 
10       $\sigma \leftarrow \sigma \cdot \exp^{c_\sigma/d} \left( \frac{\|s_\sigma\|}{\mathbb{E}\|\mathcal{N}(0, I)\|} - 1 \right)$  // update  $\sigma$ 
11       $C \leftarrow (1 - c_1 + c_h - c_\mu) C + c_1 s_c s_c^T + c_\mu \sum_{z_k \in \mathcal{P}} w_k C^{\frac{1}{2}} z_k (C^{\frac{1}{2}} z_k)^T$  // update  $C$ 
12       $x \leftarrow x + c_m \sigma C^{\frac{1}{2}} \sum_{z_k \in \mathcal{P}} w_k z_k$  // update  $x$ 
13   return  $x$ 
```

Remark: Two search paths, s_σ and s_c , are part of the algorithm state, together with x, σ , and C . They accumulate the algorithm moves across iterations.

And what are all those c_1, c_h, c_μ, \dots ?

CMA-ES Code

The full CMA-ES pseudocode:

Algorithm 8: CMA-ES

Given: $D \in \mathbb{N}_+$, $\lambda \geq 5$, $\mu \approx \lambda/2$, $w_k = w'(k) / \sum_{k=1}^{\mu} w'(k)$, $w'(k) = \log(\lambda/2 + 1/2) - \log \text{rank}(f(x_k))$, $\mu_w = 1 / \sum_{k=1}^{\mu} w_k^2$, $c_\sigma \approx \mu_w / (D + \mu_w)$, $d \approx 1 + \sqrt{\mu_w / D}$, $c_c \approx (4 + \mu_w / D) / (D + 4 + 2\mu_w / D)$, $c_1 \approx 2 / (D^2 + \mu_w)$, $c_\mu \approx \mu_w / (D^2 + \mu_w)$, $c_m = 1$.

```

1 begin
2   Initialize:  $x \in \mathbb{R}^D$ ,  $\sigma \in \mathbb{R}_+^D$ ,  $C = I$ ,  $s_\sigma = 0$ ,  $s_c = 0$ .
3   while not TerminationCondition() do
4     for  $k \in 1, \dots, \lambda$  do
5        $z_k \leftarrow \mathcal{N}(0, I)$ 
6        $x_k \leftarrow x + \sigma C^{\frac{1}{2}} \times z_k$ 
7      $\mathcal{P} \leftarrow \text{SelectBest}(\mu, \{z_k, f(x_k) \mid 1 \leq k \leq \lambda\})$ 
8      $s_\sigma \leftarrow (1 - c_\sigma)s_\sigma + \sqrt{c_\sigma(2 - c_\sigma)}\sqrt{\mu_w} \sum_{z_k \in \mathcal{P}} w_k z_k$  // search path for  $\sigma$ 
9      $s_c \leftarrow (1 - c_c)s_c + h_\sigma \sqrt{c_c(2 - c_c)}\sqrt{\mu_w} \sum_{z_k \in \mathcal{P}} w_k C^{\frac{1}{2}} z_k$  // search path for  $C$ 
10     $\sigma \leftarrow \sigma \cdot \exp^{c_\sigma/d} \left( \frac{\|s_\sigma\|}{\mathbb{E}\|\mathcal{N}(0, I)\|} - 1 \right)$  // update  $\sigma$ 
11     $C \leftarrow (1 - c_1 + c_h - c_\mu)C + c_1 s_c s_c^T + c_\mu \sum_{z_k \in \mathcal{P}} w_k C^{\frac{1}{2}} z_k (C^{\frac{1}{2}} z_k)^T$  // update  $C$ 
12     $x \leftarrow x + c_m \sigma C^{\frac{1}{2}} \sum_{z_k \in \mathcal{P}} w_k z_k$  // update  $x$ 
13  return  $x$ 
14  where  $h_\sigma = \mathbb{1}(\|s_\sigma\|^2 / D < 2 + 4/(D + 1))$ ,  $c_h = c_1(1 - h_\sigma^2)c_c(2 - c_c)$ , and  $C^{\frac{1}{2}}$  is the unique symmetric positive definite matrix obeying  $C^{\frac{1}{2}} \times C^{\frac{1}{2}} = C$ . All  $c$ -coefficients are  $\leq 1$ .

```

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 34 / 38

CMA-ES Summary

CMA-ES is **quasi parameter-free**:

- It has a lot of internal parameters, but almost all of them are carefully set by the algorithm itself.
- The user has to specify only
 - initial solution x ,
 - initial step size σ , and
 - the number of offspring λ (but even that can be set based on the search space dimension).

CMA-ES Variants:

- Reducing the local search character of CMA-ES:
 - **IPOP-CMA-ES**: Restart CMA-ES several times, making the population twice as large each time.
 - **BIPOP-CMA-ES**: Restart CMA-ES many times in 2 regimes: IPOP, and small-pop (spend similar number of evaluations in IPOP and small-pop modes).
- Reducing the number of parameters to be adapted:
 - **L-CMA-ES**: Smaller memory requirements, suitable for high-dimensional spaces, limited adaptation.
- Learning from unsuccessful mutations:
 - **Active CMA-ES**: negative weights allowed during covariance update. Gotcha: C may lose positive definiteness!

P. Pošík © 2022

A0M33EOA: Evolutionary Optimization Algorithms – 35 / 38

Relations to other algorithms

Estimation of Distribution Algorithms (**EDA**):

- CMA-ES can be considered an instance of EDA.
- EDAs template: sample from probabilistic model, and update model based on good individuals (i.e., the same as CMA-ES uses).

Natural Evolution Strategies (**NES**):

- Idea: the update of all distribution parameters should be based on the same fundamental principle.
- NES proposed as more principled alternative to CMA-ES.
- Later it was found that CMA-ES actually implements the underlying NES principle.

Information Geometric Optimization (**IGO**):

- Framework unifying many successful algorithms from discrete and continuous domains.
- CMA-ES and NES (and many EDA variants, see the next lecture) can be derived as special instances of IGO.

Summary

37 / 38

Learning outcomes

After this lecture, a student shall be able to

- perform the mapping of chromosomes from binary to real space when using binary encoding for real-parameter optimization;
- describe and exemplify the effects of such a genotype-phenotype mapping on the neighborhood structures induced by mutation and crossover;
- give examples and describe some mutation and crossover operators designed for spaces of real number vectors;
- explain the main features of ES and differences to GAs;
- explain the notation $(\mu/\rho^{\dagger}\lambda)$ -ES;
- describe the differences between mutation with isotropic, axis-parallel, and general Gaussian distribution, including the relation to the form of the covariance matrix, and the number of parameters that must be set/adapted for each of them;
- explain and use two simple methods of mutation step size adaptation (1/5 rule and self-adaptation);
- write a high-level pseudocode of CMA-ES and describe CMA-ES in the $(\mu/\rho^{\dagger}\lambda)$ notation;
- implement DE algorithm;
- explain the basic forms of DE mutation and crossover.