# Artificial Neural Networks, Symmetries and Differential Evolution

Onay Urfalioglu, Orhan Arikan Department of Electrical and Electronics Engineering, Bilkent University

October 26, 2018

#### Abstract

Neuroevolution is an active and growing research field, especially in times of increasingly parallel computing architectures. Learning methods for Artificial Neural Networks (ANN) can be divided into two groups. Neuroevolution is mainly based on Monte-Carlo techniques and belongs to the group of global search methods, whereas other methods such as backpropagation belong to the group of local search methods. ANN's comprise important symmetry properties, which can influence Monte-Carlo methods. On the other hand, local search methods are generally unaffected by these symmetries. In the literature, dealing with the symmetries is generally reported as being not effective or even yielding inferior results. In this paper, we introduce the so called Minimum Global Optimum Distance principle derived from theoretical considerations for effective symmetry breaking, applied to offline supervised learning. Using Differential Evolution (DE), which is a popular and robust evolutionary global optimization method, we experimentally show significant global search efficiency improvements by symmetry breaking.

#### 1 Introduction

Artificial Neural Networks (ANN) are general function approximators [11] and can be used to find a functional representation of a data set. Another point of view is that ANN's represent a way of data compression [2]. The compression ratio depends on the number of neurons used in the ANN which encodes the data: the less neurons at the same representation quality, the better the compression.

The estimation of the ANN-parameters is generally a computationally demanding task [23]. The corresponding Maximum-Likelihood derived cost function comprises many local optima. Therefore, local search techniques to find an optimal solution generally fail and only a suboptimal solution is found, which is a local optimum [11]. In addition, local search techniques are mainly sequential methods and parallel implementations are limited. On the other hand,

global optimization techniques based on Monte-Carlo methods such as the Genetic Algorithm (GA) [6,19], Covariance Matrix Adaptation Evolution Strategies (CMA-ES) [9,10] or Differential Evolution (DE) [20,24,28] are generally very well parallelizable. Differential Evolution is one of the most popular and robust Monte-Carlo global search methods, which outperforms many other evolutionary algorithms on a wide range of problems [3,27,29]. DE is successfully used in many engineering problems such as multiprocessor synthesis [21], optimization of radio network designs [18], training RBF networks [16], training multi layer neural networks [13] and many others.

Due to inherent symmetries in the parametric representation of ANN's, there are also multiple global optima. The multiple global optima result from point symmetries and permutation symmetries [25,26]. The effect of these symmetries on Genetic Algorithms is reported to be minimal and negligable [8]. However, there are no reports on the impact of the ANN-symmetries regarding the DE method. In this paper, we show that DE is very sensitive to multiple global optima. We derive a symmetry breaking operator based on theoretical considerations, which is optimal according to a *Minimum Global Optimum Distance* condition. In experimental studies on offline supervised learning problems, a significant improvement of up to two orders of magnitude is achieved by symmetry breaking in terms of global convergence speed. Comparisons to CMA-ES, which is a state-of-the-art evolutionary method for ANN-learning [7,22], show that CMA-ES is outperformed on complex learning problems using smaller networks which represent better compression.

## 2 Brief Review of Artificial Feedforward Neural Networks

We deal with Artificial (Feedforward) Neural Networks (ANN) for approximation of functions  $f: [-1,1]^d \to [-1,1]^q$ , having L layers (one input layer, L-2 hidden layers and one output layer) and  $N_l$  sigmoid type neurons per hidden layer l. For each neuron (l,n), we denote a parameter vector by

$$\boldsymbol{\eta}_n^l = (\boldsymbol{w}_n^l, \tau_n^l),\tag{1}$$

where  $\boldsymbol{w}_n^l$  is the weight vector and  $\tau_n^l$  is the shift scalar. The output of a tanh-type sigmoid neuron (l,n) is calculated by

$$x_n^l = \tanh(\boldsymbol{w_n^l}^\top \boldsymbol{x}^{l-1} + \tau_n^l), \tag{2}$$

where  $\mathbf{x}^l = (x_1^l, \dots, x_{N_l}^l)$  is the output vector of layer l. After all hidden layers  $l = 2, 3, \dots, L-1$  are evaluated, the final output component  $\hat{y}_n$  of the output vector  $\hat{y}$  is calculated by

$$\hat{y}_n = \mathbf{w}_n^{L^{\top}} \mathbf{x}^{L-1}, \ n = 1, .., q.$$
 (3)

We denote the parameter vector of all neurons in a layer l by  $\lambda^l$ , where

$$\boldsymbol{\lambda}^l = (\boldsymbol{\eta}_1^l, \dots, \boldsymbol{\eta}_{N_l}^l). \tag{4}$$

The total parameter vector of the whole network is given by

$$\boldsymbol{\theta}_a = (\boldsymbol{\lambda}^2, \dots, \boldsymbol{\lambda}^{L-1}, \boldsymbol{w}_1^L, \dots, \boldsymbol{w}_q^L), \tag{5}$$

where  $\boldsymbol{w}^L=(w_1^L,\dots,w_{N_L}^L)$  is the vector of the output layer weights. The function defined by the network is denoted by

$$\hat{\boldsymbol{y}} = \Omega(\boldsymbol{\theta}_a; \boldsymbol{x}), \tag{6}$$

where  $\boldsymbol{x}$  is the input vector, which equals to the output of the input layer, so that  $\boldsymbol{x}^1 \equiv \boldsymbol{x}$ .

Assuming additive normal i.i.d. noise on the available data  $(\boldsymbol{x}_k, \boldsymbol{y}_k), k = 1, ..., K$ , the ML-estimate  $\hat{\boldsymbol{\theta}}_a$  of the parameters  $\boldsymbol{\theta}_a$  is determined by the least squares solution:

$$\hat{\boldsymbol{\theta}}_a = \arg\min_{\boldsymbol{\theta}_a} \sum_{k=1}^K (\boldsymbol{y}_k - \Omega(\boldsymbol{\theta}_a; \boldsymbol{x}_k))^\top (\boldsymbol{y}_k - \Omega(\boldsymbol{\theta}_a; \boldsymbol{x}_k)).$$
 (7)

Since the output layer is linear as shown in Eqn. (3), the corresponding weights  $\boldsymbol{w}_{M,n}$  can be determined by a least squares method, as described in [17], which we adopt in this paper. This has the advantage that global search is applied only to the non-linear part of the parameter space, which generally speeds up convergence. Consequently, the parameter vector for global optimization  $\boldsymbol{\theta}$  consists of

$$\boldsymbol{\theta} = (\boldsymbol{\lambda}^2, \dots, \boldsymbol{\lambda}^{L-1}). \tag{8}$$

In the following section, we briefly review the DE method.

### 3 Brief Review of Differential Evolution

DE is one of the best general purpose evolutionary global optimization methods available. It has only linear complexity and it is known as an efficient global optimization method for continuous problem spaces. The optimization is based on a population of  $N_p$  solution candidates  $\theta_i, i \in \{1, ..., N_p\}$  where each candidate has a position in the D-dimensional search space. Initially, the solution candidates are generated randomly according to a uniform distribution within the provided intervals of the search space. The population improves by generating new positions iteratively for each candidate. For each individual  $\theta_{i,G}$ , a new trial position  $u_{i,G}$  is determined by

$$\mathbf{v}_{i,G} = \boldsymbol{\theta}_{r_1,G} + F \cdot (\boldsymbol{\theta}_{r_2,G} - \boldsymbol{\theta}_{r_3,G}) \tag{9}$$

$$\mathbf{u}_{iG} = C(\boldsymbol{\theta}_{iG}, \mathbf{v}_{iG}), \tag{10}$$

where  $r_1, r_2, r_3$  are pairwise different randomly chosen integers from the discrete set  $\{1, ..., N_p\}$  and F is a weighting scalar. The vector  $\mathbf{v}_{i,G}$  is used together with  $\boldsymbol{\theta}_{i,G}$  in the crossover operation, denoted by C(). The crossover operator copies coordinates from both  $\boldsymbol{\theta}_{i,G}$  and  $\mathbf{v}_{i,G}$  in order to create the trial vector  $\mathbf{u}_{i,G}$ . C is provided with the probability  $C_r$  to copy coordinates from  $\boldsymbol{\theta}_{i,G}$ , whereby coordinates from  $v_{i,G}$  are copied with a probability of  $1-C_r$  to  $u_{i,G}$ . Only if the new candidate  $u_{i,G}$  proves to have a lower cost then it replaces  $\theta_{i,G}$ , otherwise it is discarded.

DE includes an adaptive range scaling for the generation of solution candidates through the difference term in Equation (9). This leads to a global search with large step sizes in the case where the solution candidate vectors are widely spread within the search space due to a relatively large mean difference vector. In the case of a converging population, the mean difference vector becomes relatively small and this enables efficient fine tuning at the final phase of the optimization process. The crossover operator has a complicated role in the dynamics of the population. For example, it produces rotations that are very important when dealing with separable variables. In some cases, it can help to increase the diversity of the population or it can also speed up the convergence, depending on the problem.

### 4 Symmetries in ANN's

A symmetry is an operator  $\Phi$  which applies to the parameter vector  $\boldsymbol{\theta}$  of the ANN and leaves the output of the ANN invariant:

$$\Omega(\boldsymbol{\theta}; \boldsymbol{x}) = \Omega(\Phi(\boldsymbol{\theta}); \boldsymbol{x}) \ \forall \boldsymbol{\theta}, \boldsymbol{x}. \tag{11}$$

Non reducable ANN's comprise two types of symmetries. The first type is a *point symmetry* on the neuron parameter level, since

$$w \tanh(x) = -w \tanh(-x) \ \forall w, x. \tag{12}$$

It follows, that a point symmetry operator  $O_n^l$  defined by

$$O_n^l(\boldsymbol{\theta}) : \left\{ \begin{array}{ll} \boldsymbol{\eta}_n^l & \to & -\boldsymbol{\eta}_n^l \\ w_{i,n}^{l+1} & \to & -w_{i,n}^{l+1}, \ i = 1, \dots, N_{l+1} \end{array} \right.$$
 (13)

applied to the parameters of neuron (l,n) and the n-th weight component  $w_{i,n}^{l+1}$  to all neurons in the following layer l+1 does not change the output of the ANN. In Fig. 1, an example for the application of  $O_1^2$  is given. For each layer, the point symmetry yields  $2^{N_l}$  symmetric equivalents of the parameter vector  $\boldsymbol{\theta}$  due to the point symmetries.

The second type of symmetry is a permutation symmetry by the neuron parameters  $\eta$  and corresponding weight parameters in the next layer. A permutation operator  $P_{i,k}^l$  defined by

$$P_{j,k}^{l}(\boldsymbol{\theta}): \left\{ \begin{array}{l} \boldsymbol{\eta}_{j}^{l} \leftrightarrow \boldsymbol{\eta}_{k}^{l} \\ w_{i,j}^{l+1} \leftrightarrow w_{i,k}^{l+1}, \ i = 1, \dots, N_{l+1} \end{array} \right.$$
 (14)

leaves the output invariant. Note that  $P_{j,k}^l = P_{k,j}^l$ . In Fig. 2, the application of  $P_{1,2}^2 = P_{2,1}^2$  is illustrated. In each layer, there are  $N_l!$  symmetric equivalents

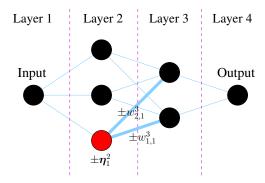


Figure 1: Application of the point symmetry operator  $O_1^2$ , which changes the signs of  $\eta_1^2$ -parameters in layer two and  $w_{n,1}^3$ -parameters in layer three, respectively.

of the parameter vector  $\boldsymbol{\theta}$  due to the permutation symmetries. The total count of symmetric equivalents per layer l is  $2^{N_l}N_l!$ . Another important property is that the length of the vector  $\boldsymbol{\theta}$  is invariant under such symmetry operators:

$$||\Phi(\boldsymbol{\theta})|| = ||\boldsymbol{\theta}|| \ \forall \Phi. \tag{15}$$

As a result, all symmetric equivalents of a global optimum lie on a hypersphere. Since such symmetries multiply the local and global optima count in the parameter space, the ultimate goal of symmetry breaking is to reduce the total number of local optima in the parameter space by avoiding all but one symmetrically equivalent space partitions. In this case, there are infinitely many ways for symmetry breaking by using the operators  $O_n^l$  and  $P_{j,k}^l$ . The differences arise from the condition on which these operators are to be applied. Fig. 3 shows different ways for breaking a point symmetry in relation to the global optimum. It can be seen that the symmetry invariant region which has maximum distances to the global optima enables the optimal separation or partitioning. This way, an optimal isolation between all symmetric equivalents of the global optima is achieved. As a result, the influence of other neighbouring global optima is decreased to a minimum, which maximizes the attraction of the global optimum of the selected partition. It can be easily shown that such a symmetry invariant region is the result of the following separation condition

$$\hat{\boldsymbol{\theta}} = \begin{cases} \boldsymbol{\theta} & \text{for } ||\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|| \le ||-\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|| \\ -\boldsymbol{\theta} & \text{for } ||\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|| > ||-\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}||, \end{cases}$$
(16)

where  $\bar{\boldsymbol{\theta}}$  denotes the selected global optimum,  $\boldsymbol{\theta}$  the parameter vector before symmetry breaking and  $\hat{\boldsymbol{\theta}}$  the parameter vector after symmetry breaking. Taking into account the additional permutation symmetry, which is independent from the point symmetry, the operator  $\Phi$  is generally composed of a chain of point symmetry and permutation symmetry operators. As an example

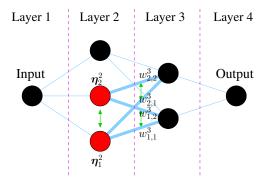


Figure 2: Application of the permutation symmetry operator  $P_{1,2}^2 = P_{2,1}^2$ , which exchanges the parameters  $\eta_1^2 \leftrightarrow \eta_2^2$  in layer two and the parameters  $w_{1,1}^3 \leftrightarrow w_{1,2}^3$ ,  $w_{2,1}^3 \leftrightarrow w_{2,2}^3$  in layer three.

 $\Phi = O_2^2 \circ P_{2,4}^3$ , applies a permutation symmetry followed by a point symmetry operator. We denote the set of all possible symmetry operators by  $\mathcal{S}$ . We generalize the separation condition (16) and define the following ideal separation, or, in other words, symmetry breaking

$$\hat{\Phi} = \arg\min_{\Phi \in \mathcal{S}} ||\Phi(\theta) - \bar{\theta}||, \ \hat{\theta} = \hat{\Phi}(\theta). \tag{17}$$

This means, the ideal separation selects the symmetry operator  $\hat{\Phi}$  from the set of all possible symmetry operators which minimizes the distance of the parameter vector  $\hat{\boldsymbol{\theta}}$  to a global optimum  $\bar{\boldsymbol{\theta}}$ .

### 4.1 Approximation of the ideal separation

There are two problems for the practical application of the ideal separation. First, the global optimum is not known a priori. However, iterative algorithms like DE produce intermediate results at each iteration, which can be regarded as an approximation of the global optimum. This approximation becomes better with increasing iteration number. Therefore, we propose to choose the best individual of the DE-population  $\tilde{\theta}$  at each iteration as an approximation for the global optimum. Second, the brute force method for finding an optimal solution to (17) has exponential complexity. Instead of trying to find the optimal symmetry breaking operator, we propose to apply only one single symmetry operator  $O_n^l$  or  $P_{j,k}^l$  at a time. The symmetry operator, layer and neuron are randomly selected. The proposed heuristic only approximates  $\hat{\Phi}$ , but this approximation improves with the iteration count. For each neuron (l,n), we define a symmetry relevant parameter block  $\beta_n^l$  as

$$\boldsymbol{\beta}_{n}^{l} = (\boldsymbol{\eta}_{n}^{l}, w_{1,n}^{l+1}, \dots, w_{N_{l+1},n}^{l+1}), \ l = 2, \dots, L-2 \text{ and } \boldsymbol{\beta}_{n}^{L-1} = (\boldsymbol{\eta}_{n}^{L-1}).$$
 (18)

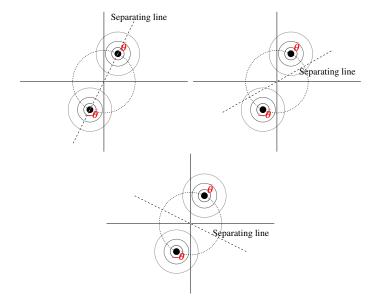


Figure 3: Example for a point symmetry in 2-D, where  $f(\theta) = f(-\theta) \ \forall \theta$ . The plots show worst case (left), suboptimal (middle) and optimal separation lines (right) for point symmetry breaking. The separating line divides the parameter space in two parts, where each partition contains a global optimum  $(\bar{\theta} \text{ and } -\bar{\theta})$ .

Given a parameter vector  $\boldsymbol{\theta}$  and an approximation of the global optimum  $\boldsymbol{\dot{\theta}}$  with corresponding parameter blocks  $\tilde{\beta}_n^l$ , the pseudocode (1) describes the proposed heuristic. A property of this heuristic is that the symmetry is not completely and not uniquely broken and the resulting modified parameters may belong to different partitions over time. As also  $\boldsymbol{\tilde{\theta}}$  changes over time, the final convergence result will be a random global optimum among all other possible global optima.

The proposed symmetry breaking method is always applied on each individual's position  $\theta_{i,G}$  (see Eqn. (10)) at each iteration prior to the application of Eqn. (10).

# 5 Experiments

In this section, we introduce results of experiments to demonstrate the performance improvements by symmetry breaking. The following methods are compared in regression and classification tests: Differential Evolution (DE), Covariance Matrix Adaptation Evolution Strategies (CMA-ES) and DE with symmetry breaking (DE-SB). With a D-dimensional parameter space, all tests are performed with following settings:

• DE, DE-SB settings: F = 0.5,  $C_r = 0.9$ , initial population is randomly

Algorithm 1 Proposed heuristic for breaking symmetry. Parameters  $\mu, l, n, m$  are sampled from corresponding discrete uniform distributions U(.). A symmetry operator X is only applied to the parameter vector  $\boldsymbol{\theta}$  when it decreases the distance to the global optimum  $\tilde{\boldsymbol{\theta}}$ , i.e.,  $||X(\boldsymbol{\theta}) - \tilde{\boldsymbol{\theta}}|| < ||\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}||$ . Algorithm input:  $\boldsymbol{\theta}$  and  $\tilde{\boldsymbol{\theta}}$ , effect: (eventually) modify parameter vector  $\boldsymbol{\theta}$ .

```
symmetry operator selection: sample \mu \sim U(\{0,1\})
hidden layer selection: sample l = U(\{2, ..., L-1\})
neuron selection: sample n = U(\{1, ..., N_l\}) in hidden layer l
if \mu = 0 then
   [point symmetry operator selected]
   // would the point symmetry operator O_n^l decrease the distance? (||O_n^l(\boldsymbol{\theta}) - \tilde{\boldsymbol{\theta}}|| < 1)
   calculate distance-square for NOT applying the operator O_n^l: D_1 = ||\boldsymbol{\beta}_n^l - \tilde{\boldsymbol{\beta}}_n^l||^2
   calculate distance-square for applying the operator O_n^l: D_2 = ||-\beta_n^l - \tilde{\beta}_n^l||^2
   if D_1 > D_2 then
      apply point symmetry operator O_n^l: set \beta_n^l = -\beta_n^l
   end if
else
   [permutation symmetry operator selected]
   second neuron selection: sample m = U(\{1, ..., N_l\}) in hidden layer l
   // would the permutation operator P_{m,n}^l decrease the distance? (||P_{m,n}^l(\theta) - \tilde{\theta}|| < 1)
   calculate distance-square for NOT applying the operator P_{m,n}^l: D_1 = ||\boldsymbol{\beta}_n^l||
   |\hat{\beta}_{n}^{l}||^{2} + ||\hat{\beta}_{m}^{l} - \tilde{\beta}_{m}^{l}||^{2}
   calculate distance-square for applying the operator P_{m,n}^l: D_2 = ||\boldsymbol{\beta}_n^l - \tilde{\boldsymbol{\beta}}_m^l||^2 +
   ||\boldsymbol{\beta}_m^l - \tilde{\boldsymbol{\beta}}_n^l||^2
   if D_1 > D_2 then
      apply permutation symmetry operator P_{m,n}^l: swap \beta_m^l \leftrightarrow \beta_n^l
   end if
end if
```

generated in D-dim. hypercube  $[-1,1]^D$  (uniformly),

• CMA-ES settings: we used suggested settings for enhanced global search abilities, mentioned in the C-code reference implementation.

Given a parameter  $\boldsymbol{\theta}$  and a data set  $(\boldsymbol{x}_i, \boldsymbol{y}_i)$ , the ANN produces a cost  $\epsilon$  defined by the Mean Squared Error (MSE), which is derived from Eqn. (7):

$$\epsilon = \frac{1}{K \cdot q} \sum_{k=1}^{K} (\boldsymbol{y}_k - \Omega(\boldsymbol{\theta}; \boldsymbol{x}_k))^{\top} (\boldsymbol{y}_k - \Omega(\boldsymbol{\theta}; \boldsymbol{x}_k)).$$
 (19)

In order to limit the D-dimensional parameter space to a feasible region, we apply a penalty approach. Due to the length-invariance by the symmetry operators as shown in Eqn. (15), the feasible region is defined by a hypersphere. In case of  $||\boldsymbol{\theta}|| > \sqrt{D}$ , the cost function (19) is evaluated at a rescaled parameter vector  $\frac{\boldsymbol{\theta}}{||\boldsymbol{\theta}||}$  and a penalty term  $50(||\boldsymbol{\theta}|| - \sqrt{D})$  is added to the cost  $\epsilon$ .

#### 5.1 Regression problems

All utilized target functions are defined over  $\boldsymbol{x} \in [-1,1]^d$  and map to [-1,1]. A data set, which consists of training and test samples, is generated by sampling  $\boldsymbol{x}_i$  from a d-dimensional uniform distribution  $U^d(-1,1)$ 

$$\boldsymbol{x}_i \sim U^d(-1,1),\tag{20}$$

and adding normal distributed noise with zero mean and variance  $\sigma^2$  to the function values  $y_i$ 

$$y_i = f(\mathbf{x}_i) + \mu, \ \mu \sim \mathcal{N}(0, \sigma^2), \ \sigma = 5 \times 10^{-3}.$$
 (21)

All data sets of regression problems contain 200 training samples and 200 test samples. Tab. 1 shows the utilized functions and corresponding data sets for regression tests. The population size  $N_p$  depends on the problem and is manually adapted accordingly. For each problem, the global optimization is applied in 50 independent runs and the mean required number of ANN-evaluations (MFE) and corresponding standard deviation  $\sigma_{\rm MFE}$  to reach the error threshold  $\epsilon_0$  are determined. The population size manually is adapted such that no one of the 50 runs fails to reach the error threshold and the MFE is kept minimal. We declare a global optimum as found by reaching the error threshold. We define a robustness  $\rho$  as

$$\rho = \frac{\text{number of successfull runs}}{\text{total number of runs}}.$$
 (22)

Results for the most complex settings (least number of hidden neurons) are shown in Tab. 2. Fig. 4 shows results for the required number of ANN-evaluations (MFE) over the number of hidden neurons used in the ANN.

#### 5.2 Classification problems

The following classification problems are used: **iris** data set [4,5], **tic-tac-toe** data set [1], **balance** [12,14] and a more challanging problem defined by the **two-spirals** [15] data set. Samples are divided into a training set and a test set following the format a/b. The selection process iteratively and sequentially puts a samples into the training set and the following b samples into the test set, until no samples are available. A winner-takes-all scheme is applied to

Table 1: Target funtions, corresponding data sets, error thresholds and typical test set MSE's of function regression experiments. There was no significant variation of the MSE results by the choice of the learning method.

function	data set	training / test samples	$\epsilon_0$	test set MSE
$(x+0.5)^2(0.1+(x+0.65)^2)$	syn5	200 / 200	$5 \times 10^{-5}$	$6.27 \times 10^{-5} \pm 3.2 \times 10^{-6}$
$\sin(10x)/(10x)$	sinc	200 / 200	$5 \times 10^{-5}$	$5.99 \times 10^{-5} \pm 7.6 \times 10^{-6}$
$x/2 + \sin(10x)/(10x)$	incsinc	200 / 200	$5 \times 10^{-5}$	$5.88 \times 10^{-5} \pm 7.1 \times 10^{-6}$
$\sin(5r)/(15r), r = \sqrt{x^2 + y^2}$	sinc2d	200 / 200	$5 \times 10^{-5}$	$5.74 \times 10^{-5} \pm 6.7 \times 10^{-6}$

Table 2: Results of Mean ANN (Function) Evaluations (MFE) from 50 independent runs and population sizes of regression tests. Better values are highlighted in boldface. Note that on sinc (1-5-1) and incsinc (1-5-1), CMA-ES failed to find the global optimum on some runs, even though a large population size of 10000 was used. On sinc2d (2-3-1-3-1), CMA-ES failed in all 50 runs.

		DE	DE-SB	CMA-ES	
data set	topology	[population size, robustness] MFE	[population size, robustness] MFE	[population size, robustness] MFE	$\frac{\text{MFE (DE)}}{\text{MFE (DE-SB)}}$
syn5	1-3-1	$[80, 1]7.30 \times 10^4 \pm 1.8 \times 10^4$	$[80, 1]3.20 \times 10^4 \pm 9.4 \times 10^3$	[160, 1] <b>1.17×10<sup>4</sup> ± 2.6×10<sup>3</sup></b>	2.3
sinc	1-5-1	$[1400, 1]9.22 \times 10^7 \pm 3.1 \times 10^7$	$[160, 1]4.47 \times 10^5 \pm 2.8 \times 10^5$	$[10^4, 0.36]3.78 \times 10^6 \pm 3.7 \times 10^5$	206
sinc	1-6-1	$[800, 1]4.19 \times 10^7 \pm 2.2 \times 10^7$	$[60, 1]1.52 \times 10^5 \pm 4.5 \times 10^4$	[100, 1] <b>2.67</b> × <b>10</b> <sup>4</sup> ± <b>5.3</b> × <b>10</b> <sup>3</sup>	276
incsinc	1-5-1	$[1600, 1]1.36 \times 10^8 \pm 6.0 \times 10^7$	$[200, 1]8.04 \times 10^{5} \pm 1.8 \times 10^{5}$	$[10^4, 0.8]7.76 \times 10^6 \pm 1.9 \times 10^6$	169
incsinc	1-6-1	$[800, 1]4.36 \times 10^7 \pm 2.6 \times 10^7$	$[56, 1]1.18 \times 10^5 \pm 5.0 \times 10^4$	[100, 1]2.86×10 <sup>4</sup> ± 8.2×10 <sup>3</sup>	370
sinc2d	2-3-1-3-1	$[120, 1]1.93 \times 10^5 + 4.2 \times 10^4$	$[120, 1]$ 1, $22 \times 10^5 + 2.4 \times 10^4$	[10 <sup>4</sup> , 0] NA	1.58

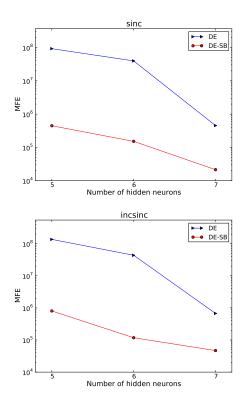


Figure 4: Required mean function evaluations (MFEs) to find the global optimum depending on the number of hidden neurons  $N_2$ . The corresponding network topology is  $(1 - N_2 - 1)$ .

distinguish different classes. The output vector  $\boldsymbol{y}$  of a sample designating class i has the following format

$$y_j = \begin{cases} 1 & \text{for } j = i \\ 0 & \text{else.} \end{cases}$$
 (23)

As in the case of function regression experiments, on each classification problem, a predefined error threshold  $\epsilon_0$  is used as a termination criterion for the learning process. Tab. 3 shows the error thresholds and typical classification sucess rates. Tab. 4 shows the results on the classification data sets. Comparing DE

Table 3: Data sets, number of training and test samples, the sample selection format, corresponding error thresholds and typical classification sucess rates of classification experiments.

data set	training / test samples	sample selection	$\epsilon_0$	test set classification success (%)
iris	75/75	1/1	0.011	$99.4 \pm 0.6$
tic-tac-toe	479/479	1/1	0.08	$92.1 \pm 0.8$
balance	312/313	1/1	0.001	$100 \pm 0$
two-spirals	97/97	2/2	0.07	$91.2 \pm 1.4$

Table 4: Results and population sizes of classification tests. Better values are highlighted in boldface. Note that on balance, CMA-ES failed to find the global optimum on some runs, even though a large population size of 10000 was used. On two-spirals (2-10-1-10-2) and two-spirals (2-9-1-9-2), CMA-ES failed in all 50 runs.

		DE	DE-SB	CMA-ES	
data set	topology	[population size, robustness] MFE	[population size, robustness] MFE	[population size, robustness] MFE	$\frac{\text{MFE (DE)}}{\text{MFE (DE-SB)}}$
iris	4-3-3	$[40, 1]1.46 \times 10^4 \pm 3.3 \times 10^3$	$[40, 1]1.24 \times 10^4 \pm 2.5 \times 10^3$	[20, 1]3.74×10 <sup>3</sup> ± 4.2×10 <sup>2</sup>	1.18
tic-tac-toe	9-8-2	$[230, 1]2.70 \times 10^6 \pm 4.7 \times 10^5$	$[160, 1]9.74 \times 10^5 \pm 2.2 \times 10^4$	[800, 1] <b>3.29</b> ×10 <sup>5</sup> ± <b>3.9</b> ×10 <sup>4</sup>	2.78
balance	4-5-1-5-3	$[520, 1]1.38 \times 10^7 \pm 9.2 \times 10^6$	[160, 1]8.06×10 <sup>5</sup> ± 3.2×10 <sup>5</sup>	$[10^4, 0.78]5.67 \times 10^6 \pm 5.6 \times 10^6$	17.1
two-spirals	2-10-1-10-2	$[200, 1]2.24 \times 10^{7} \pm 4.8 \times 10^{6}$	$[120, 1]6.74 \times 10^{6} \pm 2.5 \times 10^{6}$	[10 <sup>4</sup> , 0] NA	3.32
two-spirals	2-9-1-9-2	$[240, 1]4.26 \times 10^7 \pm 1.5 \times 10^7$	$[120, 1]8.40 \times 10^6 \pm 2.6 \times 10^6$	$[10^4, 0] \text{ NA}$	5.07

and DE-SB, symmetry breaking consistently improves global search efficiency in all experiments. On reducable networks with a larger number of hidden neurons, as in **sinc** and **incsinc**, CMA-ES is superior and capable of robustly finding a solution. However, on networks with the smallest number of hidden neurons representing best compression, DE-SB outperforms CMA-ES. Furthermore, on complex problems with deeper network topologies, CMA-ES seems to have difficulties to robustly find the global optimum, even with a very large population size. On these type of problems, the true global search character of DE pays off.

#### 6 Conclusions

It is shown that symmetries in ANN-parameter space do affect the performance of Differential Evolution (DE). From theoretical considerations, we derive an ideal operator for breaking these symmetries. This ideal operator requires knowledge about the global optimum of the parameter space. Since the global

optimum is not known a priori, the ideal operator is not applicable. Another concern is that a brute force implementation of the ideal operator has exponential complexity. Therefore, we propose a heuristic to approximate the ideal operator, which has negligable overhead. Unlike the CMA-ES method, which has at least quadratic complexity, the proposed DE with symmetry breaking (DE-SB) has linear complexity and is generally applicable on very high dimensional parameter spaces.

Experimental studies on a priori fixed topology networks indicate a significant improvement over standard DE in terms of required mean number of ANN-evaluations. Compared to CMA-ES, which is a state-of-the-art method, we achieve superior results especially on complex problems and smaller networks, which represent a better compression at comparable approximation quality. We believe that other global optimization methods may also significantly benefit from symmetry breaking. Specifics of how each method should realize a symmetry breaking heuristic requires further research.

#### References

- [1] D. W. Aha. UCI machine learning repository, 1991.
- [2] M. A. Arbib, editor. The Handbook of Brain Theory and Neural Networks. MIT Press, Cambridge, MA, USA, 2002.
- [3] H. Bersini, M. Dorigo, S. Langerman, G. Seront, and L. M. Gambardella. Results of the first international contest on evolutionary optimisation (1st iceo). In *International Conference on Evolutionary Computation*, pages 611–615, 1996.
- [4] B. V. Dasarathy. Nosing around the neighborhood: A new system structure and classification rule for recognition in partially exposed environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence, PAMI*, 2:67–71, 1980.
- [5] R. A. Fisher and M. Marshall. UCI machine learning repository, 1988.
- [6] D. E. Goldberg. Genetic Algorithms in Search, Optimization, and Machine Learning. Addison-Wesley Professional, 1 edition, January 1989.
- [7] F. Gomez, J. Schmidhuber, and R. Miikkulainen. Accelerated neural evolution through cooperatively coevolved synapses. J. Mach. Learn. Res., 9:937–965, 2008.
- [8] S. Haflidason and R. Neville. On the significance of the permutation problem in neuroevolution. In GECCO '09: Proceedings of the 11th Annual conference on Genetic and evolutionary computation, pages 787–794, New York, NY, USA, 2009. ACM.

- [9] N. Hansen, S. D. Müller, and P. Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evol. Comput.*, 11(1):1–18, 2003.
- [10] N. Hansen and A. Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: the covariance matrix adaptation. In *Proc. of the 1996 IEEE Int. Conf. on Evolutionary Computation*, pages 312–317, Piscataway, NJ, 1996. IEEE Service Center.
- [11] S. Haykin. Neural Networks: A Comprehensive Foundation. Prentice Hall, 2nd edition, July 1998.
- [12] T. Hume. UCI machine learning repository, 1994.
- [13] J. Ilonen, J.-K. Kamarainen, and J. Lampinen. Differential evolution training algorithm for feed-forward neural networks. *Neural Process. Lett.*, 17(1):93–105, 2003.
- [14] D. Klahr and R. S. Siegler. The representation of children's knowledge. In *Advances in child development and behavior* (, pages 61–116, 1978.
- [15] K. J. Lang and M. J. Witbrock. Learning to tell two spirals apart. In Proceedings 1988 Connectionist Models Summer School, pages 52–59, Los Altos, CA, 1988. Morgan Kaufmann.
- [16] J. Liu, J. Mattila, and J. Lampinen. Training rbf networks using a de algorithm with adaptive control. *Tools with Artificial Intelligence, IEEE International Conference on*, 0:673–676, 2005.
- [17] T. Masters. Practical Neural Networks Recipes in C++. Academic Press, 1993.
- [18] S. P. Mendes, J. A. G. Pulido, M. A. V. Rodriguez, M. D. J. Simon, and J. M. S. Perez. A differential evolution based algorithm to optimize the radio network design problem. In *E-SCIENCE '06: Proceedings of the Second IEEE International Conference on e-Science and Grid Computing*, page 119, Washington, DC, USA, 2006. IEEE Computer Society.
- [19] Z. Michalewicz. Genetic algorithms + data structures = evolution programs (2nd, extended ed.). Springer-Verlag New York, Inc., New York, NY, USA, 1994.
- [20] K. V. Price. Differential evolution: a fast and simple numerical optimizer. In Biennial Conference of the North American Fuzzy Information Processing Society, NAFIPS, pages 524–527. IEEE Press, New York. ISBN: 0-7803-3225-3, June 1996.
- [21] A. Rae and S. Parameswaran. Application-specific heterogeneous multiprocessor synthesis using differential-evolution. In *ISSS '98: Proceedings of the 11th international symposium on System synthesis*, pages 83–88, Washington, DC, USA, 1998. IEEE Computer Society.

- [22] N. T. Siebel, J. Boetel, and G. Sommer. Efficient neural network pruning during neuro-evolution. In *Proceedings of 2009 International Joint Conference on Neural Networks (IJCNN 2009)*, Atlanta, USA, pages 2920–2927, June 2009.
- [23] J. Síma. Minimizing the quadratic training error of a sigmoid neuron is hard. In ALT '01: Proceedings of the 12th International Conference on Algorithmic Learning Theory, pages 92–105, London, UK, 2001. Springer-Verlag.
- [24] R. Storn and K. Price. Differential evolution a simple and efficient adaptive scheme for global optimization over continuous spaces. *Technical Report TR-95-012, ICSI*, Mar. 1995.
- [25] H. J. Sussmann. Uniqueness of the weights for minimal feedforward nets with a given input-output map. *Neural Netw.*, 5(4):589–593, 1992.
- [26] D. Thierens. Non-redundant genetic coding of neural networks. In In Proceedings of the 1996 IEEE International Conference on Evolutionary Computation, pages 571–575. IEEE Press, 1996.
- [27] T. Tušar and B. Filipič. Differential evolution versus genetic algorithms in multiobjective optimization. In *EMO'07: Proceedings of the 4th international conference on Evolutionary multi-criterion optimization*, pages 257–271, Berlin, Heidelberg, 2007. Springer-Verlag.
- [28] J. Vesterstrom and R. Thomsen. A comparative study of differential evolution, particle swarm optimization, and evolutionary algorithms on numerical benchmark problems. In *Evolutionary Computation*, 2004. CEC2004. Congress on, volume 2, pages 1980–1987 Vol.2, June 2004.
- [29] X. Xu and Y. Li. Comparison between particle swarm optimization, differential evolution and multi-parents crossover. Computational Intelligence and Security, International Conference on, 0:124–127, 2007.