

Université Libre de Bruxelles

*Institut de Recherches Interdisciplinaires
et de Développements en Intelligence Artificielle*

**Training feed-forward neural networks
with ant colony optimization:
An application to pattern classification**

Christian BLUM and Krzysztof SOCHA

IRIDIA – Technical Report Series

Technical Report No.
TR/IRIDIA/2005-038

December 2005

IRIDIA – Technical Report Series
ISSN 1781-3794

Published by:

IRIDIA, *Institut de Recherches Interdisciplinaires
et de Développements en Intelligence Artificielle*
UNIVERSITÉ LIBRE DE BRUXELLES
Av F. D. Roosevelt 50, CP 194/6
1050 Bruxelles, Belgium

Technical report number TR/IRIDIA/2005-038

The information provided is the sole responsibility of the authors and does not necessarily reflect the opinion of the members of IRIDIA. The authors take full responsibility for any copyright breaches that may result from publication of this paper in the IRIDIA – Technical Report Series. IRIDIA is not responsible for any use that might be made of data appearing in this publication.

Training feed-forward neural networks with ant colony optimization: An application to pattern classification

Christian Blum
ALBCOM, LSI
Universitat Politècnica de Catalunya
Barcelona, Spain
Email: cblum@lsi.upc.edu

Krzysztof Socha
IRIDIA
Université Libre de Bruxelles
Brussels, Belgium
Email: ksocha@ulb.ac.be

Abstract—Ant colony optimization is an optimization technique that was inspired by the foraging behaviour of real ant colonies. Originally, the method was introduced for the application to discrete optimization problems. Recent research efforts led to the development of algorithms that are also applicable to continuous optimization problems. In this work we present one of the most successful variants for continuous optimization and apply it to the training of feed-forward neural networks for pattern classification. For evaluating our algorithm we apply it to classification problems from the medical field. The results show, first, that our algorithm is comparable to specialized algorithms for neural network training, and second, that our algorithm has advantages over other general purpose optimizers.

I. INTRODUCTION

Pattern classification is an important real-world problem. In the medical field, for example, pattern classification problems arise when physicians are interested in reliable classifiers for diseases based on a number of measurements. Feed-forward neural networks (NNs) are commonly used systems for the task of pattern classification [5], but require prior configuration. Generally the configuration problem consists hereby of two parts: First, the structure of the feed-forward NN has to be determined. Second, the numerical weights of the neuron connections have to be determined such that the resulting classifier is as correct as possible. In this work we focus only on the second part, namely the optimization of the connection weights. We adopt the NN structures from earlier works on the same subject.

Ant colony optimization (ACO) is an optimization technique that was introduced for the application to discrete optimization problem in the early 90's by M. Dorigo and colleagues [9], [10], [11]. The origins of ant colony optimization are in a field called swarm intelligence (SI) [6], which studies the use of certain properties of social insects, flocks of birds, or fish schools, for tasks such as optimization. The inspiring source of ACO is the foraging behaviour of real ant colonies. When searching for food, ants initially explore the area surrounding their nest in a random manner. While moving, ants leave a chemical pheromone trail on the ground.

As soon as an ant finds a food source, it evaluates the quantity and the quality of the food and carries some of it back to the nest. During the return trip, the quantity of pheromone that an ant leaves on the ground may depend on the quantity and quality of the food. The pheromone trails guide other ants to the food source. It has been shown in [8] that the indirect communication between the ants via pheromone trails enables them to find shortest paths between their nest and food sources. The shortest path finding capabilities of real ant colonies are exploited in artificial ant colonies for solving optimization problems.

While ACO algorithms were originally introduced to solve discrete optimization (i.e., combinatorial) problems, their adaptation to solve continuous optimization problems enjoys an increasing attention. Early applications of the ants metaphor to continuous optimization include algorithms such as Continuous ACO (CACO) [2], the API algorithm [17], and Continuous Interacting Ant Colony (CIAC) [12]. However, all these approaches follow rather loosely the original ACO framework. The latest approach, which is at the same time the approach that is closest to the spirit of ACO for combinatorial problems, was proposed in [20]. In this work we extend this approach, and apply it to the problem of optimizing the weights of feed-forward NNs for the task of pattern classification.

The outline of our work is as follows. In Section 2 we shortly present the structure of feed-forward NNs for the purpose of pattern classification. Then, in Section 3 we present the ACO algorithm, while in Section 4 we compare our algorithm to methods specialized for feed-forward NN training, as well as to a genetic algorithm. Finally, in Section 5 we offer a conclusion and a glimpse of future work.

II. FEED-FORWARD NEURAL NETWORKS FOR PATTERN CLASSIFICATION

A dataset for pattern classification consists of a number of patterns together with their correct classification. Each pattern

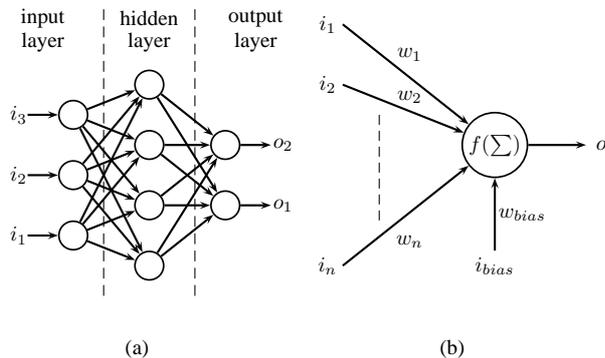


Fig. 1. (a) shows a feed-forward NN with one hidden layer of neurons. Note that each neuron of a certain layer is connected to each neuron of the next layer. (b) shows one single neuron (from either the hidden layer, or the output layer). The neuron receives inputs (i.e., signals i_l , weighted by weights w_l) from each neuron of the previous layer. Additionally, it receives a so-called bias input i_{bias} with weight w_{bias} . The transfer function $f(\Sigma)$ of a neuron transforms the sum of all the weighted inputs into an output signal, which serves as input for all the neurons of the following layer. Input signals, output signals, biases and weights are real values.

consists of a number of measurements (i.e., numerical values). The goal consists in generating a classifier that takes the measurements of a pattern as input, and provides its correct classification as output. A popular type of classifier are feed-forward neural networks (NNs).

A feed-forward NN consists of an input layer of neurons, an arbitrary number of hidden layers, and an output layer (for an example, see Figure 1). Feed-forward NNs for pattern classification purposes consist of as many input neurons as the patterns of the data set have measurements, i.e., for each measurement there exists exactly one input neuron. The output layer consists of as many neurons as the data set has classes, i.e., if the patterns of a medical data set belong to either the class normal or to the class pathological, the output layer consists of two neurons. Given the weights of all the neuron connections, in order to classify a pattern, one provides its measurements as input to the input neurons, propagates the output signals from layer to layer until the output signals of the output neurons are obtained. Each output neuron is identified with one of the possible classes. The output neuron that produces the highest output signal classifies the respective pattern.

The process of generating a NN classifier consists of determining the weights of the connections between the neurons such that the NN classifier shows a high performance. Since the weights are real-valued, this is a continuous optimization problem of the following form: Given are n decision variables $\{X_1, \dots, X_n\}$ with continuous domains. These domains are not restricted, i.e., each real number is feasible. Furthermore, the problem is unconstrained, which means that the variable settings do not depend on each other. Sought is a solution that minimizes the objective function called *square error percentage (SEP)*:

$$SEP = 100 \frac{O_{\max} - O_{\min}}{n_0 n_p} \sum_{p=1}^{n_p} \sum_{i=1}^{n_0} (t_i^p - o_i^p)^2, \quad (1)$$

where O_{\max} and O_{\min} are respectively the maximum and minimum values of the output signals of the output neurons, n_p represents the number of patterns, n_0 is the number of output neurons, and t_i^p and o_i^p represent respectively the expected and actual values of output neuron i for pattern p .

III. ANT COLONY OPTIMIZATION FOR CONTINUOUS OPTIMIZATION

ACO algorithms are iterative methods that try to solve optimization problems as follows. At each iteration candidate solutions are probabilistically constructed by sampling a probability distribution over the search space. Then, this probability distribution is modified using the better ones among the constructed solutions. The goal is to bias over time the sampling of solutions to areas of the search space that contain high quality solutions.

In ACO algorithms for discrete optimization problems, the probability distribution is discrete and is derived from artificial pheromone information. In a way, the pheromone information represents the stored search experience of the algorithm. In contrast, our ACO algorithm for continuous optimization, henceforth denoted by ACO*, utilizes a continuous probability density function (PDF). This density function is – for each solution construction – produced from a population P of solutions that the algorithm keeps at all times. The management of this population works as follows. Before the start of the algorithm, the population—whose size k is a parameter of the algorithm—is filled with random solutions. Even though the domains of the decision variables are not restricted, we used the initial interval $[-1, 1]$ for the sake of simplicity. Then, at each iteration a set of m solutions is generated and added to P . The same number of the worst solutions are removed from P . This biases the search process towards the best solutions found during the search.

For constructing a solution an ant acts as follows. First, it transforms the original set of decision variables $\mathbf{X} = \{X_1, \dots, X_n\}$ into a set of temporary variables $\mathbf{Z} = \{Z_1, \dots, Z_n\}$. The purpose of introducing temporary variables is to improve the algorithm's performance by limiting the correlation between decision variables. Note that this transformation also affects the population of solutions: All the solutions are transformed to the new coordinate system as well. The method of transforming the set of decision variables is presented towards the end of this section.

Then, at each construction step $i = 1, \dots, n$, the ant chooses a value for decision variable Z_i . For performing this choice it uses a Gaussian kernel PDF, which is a weighted superposition of several Gaussian functions. For a decision variable Z_i the Gaussian kernel G_i is given as follows:

$$G_i(z) = \sum_{j=1}^k \omega_j \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(z-\mu_j)^2}{2\sigma_j^2}}, \quad \forall z \in \mathbf{R}, \quad (2)$$

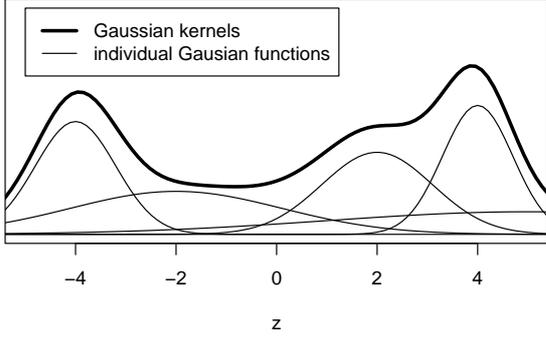


Fig. 2. An example of a Gaussian kernels PDF consisting of five separate Gaussian functions.

where the j -th Gaussian function is derived from the j -th member of the population P (remember that k is the size of P). Note that $\vec{\omega}$, $\vec{\mu}$, and $\vec{\sigma}$ are vectors of size k . Hereby, $\vec{\omega}$ is the vector of weights, whereas $\vec{\mu}$ and $\vec{\sigma}$ are the vectors of means and standard deviations respectively. Figure 2 presents an example of a Gaussian kernel PDF consisting of five separate Gaussian functions.

Sampling directly the Gaussian kernel PDF as defined in Equation 2 is problematic. Therefore we accomplish it with the following procedure. It may be proven that this procedure is exactly equivalent to sampling the PDF G_i directly.

Before starting a solution construction, we choose exactly one of the Gaussian functions j , which is then used for all n construction steps.¹ A Gaussian function j^* is chosen with the following probability distribution:

$$\mathbf{p}_j = \frac{\omega_j}{\sum_{l=1}^k \omega_l}, \forall j = 1, \dots, k, \quad (3)$$

where ω_j is the weight of Gaussian function j , which is obtained as follows. All solutions in P are ranked according to their quality (i.e., their SEP value) with the best solution having rank 1. Assuming the rank of the j -th solution in P to be r , the weight ω_j of the j -th Gaussian function is calculated according to the following formula:

$$\omega_j = \frac{1}{qk\sqrt{2\pi}} \cdot e^{-\frac{(r-1)^2}{2q^2k^2}}, \quad (4)$$

which essentially defines the weight to be a value of the Gaussian function with the argument of rank r , with mean in 1.0 and standard deviation of qk , where q is a parameter of the algorithm. When parameter q is small, the best-ranked solutions are strongly preferred, and when it is larger, the probability becomes more uniform. Thanks to using the ranks instead of the actual fitness function values, the algorithm is not sensitive to the scaling of the fitness function.

The sampling of the chosen Gaussian function j^* may be done using a random number generator that is able to

¹Note that this has also the advantage that it allows to exploit the (possibly existing) correlation between the variables.

generate random numbers according to a parameterized normal distribution, or by using a uniform random generator in conjunction with (for instance) the Box-Muller method [7]. However, before doing that we have to specify the mean and the standard deviation of the j^* -th Gaussian function. As mean μ_{j^*} we choose the value of the i -th decision variable in solution j^* . It remains to specify the standard deviation σ_{j^*} . For doing that we calculate the average distance of the other population members from the j^* -th solution and multiply it by a parameter ρ , which regulates the speed of convergence:

$$\sigma_{j^*} = \rho \sum_{l=1}^k \sqrt{(z_i^l - z_i^{j^*})^2} \quad (5)$$

Parameter ρ has a role similar to the pheromone evaporation rate ρ in the combinatorial ACO. The higher the value of $\rho \in (0, 1)$, the lower the convergence speed of the algorithm, and hence the lower the learning rate. Since this whole process is done for each dimension (i.e., decision variable) in turn, each time the distance is calculated only with the use of one single dimension (the rest of them are discarded). This ensures that the algorithm is able to adapt to convergence, but also allows the handling of problems that are scaled differently in different directions.

Finally it remains to explain how the set of temporary decision variables \mathbf{Z} is created from the original set \mathbf{X} .² An obvious choice for adapting the coordinate system to the distribution of population P would be the Principal Component Analysis (PCA) [14]. Although PCA works very well for reasonably regular distributions, its performance is no longer that interesting in case of more complex functions. The mechanism that we designed instead, is relatively simple. Each ant at each step of the construction process chooses a direction. The direction is chosen by randomly selecting a solution u from the population that is reasonably far away from the solution j^* chosen as mean of the PDF. Then, the vector j^*u becomes the chosen direction. The probability of choosing solution u (having solution j^* chosen as the mean of the PDF) is the following:

$$\mathbf{p}(u|j^*) = \frac{d(u, j^*)^4}{\sum_{l=1}^k d(l, j^*)^4}, \quad (6)$$

where function $d(\cdot)$ is the function that returns the distance between two members of the population P . Once this vector is chosen, the new orthogonal basis for the ant's coordinate system is created using the Gram-Schmidt process [13]. It takes as input all the (already orthogonal) directions chosen in earlier ant's steps and the newly chosen vector. The remaining missing vectors (for the remaining dimensions) are chosen randomly. Then, all the current coordinates of all the solutions in the population are rotated and recalculated according to this

²Note that ACO algorithms in general do not exploit correlation information between different decision variables (or components). In ACO*, due to the specific way the search experience is stored (i.e., as a population of solutions), it is in fact possible to take into account the correlation between the decision variables.

TABLE I

SUMMARY OF THE NN STRUCTURES THAT WE USE FOR THE THREE DATA SETS. IN THE LAST TABLE COLUMN IS GIVEN THE NUMBER OF WEIGHTS TO BE OPTIMIZED FOR EACH TACKLED PROBLEM. NOTE THAT FOR THE CALCULATION OF THIS NUMBER, THE BIAS INPUTS OF THE NEURONS HAVE ALSO TO BE TAKEN INTO ACCOUNT.

Data set	Inp. layer	Hid. layer	Outp. layer	# of weights
Cancer1	9	6	2	74
Diabetes1	8	6	2	68
Heart1	35	6	2	230

new orthogonal base resulting in the set of new temporary variables \mathbf{Z} . Only then is the ant able to measure the average distance, and subsequently to sample from the PDF (as it can now calculate the mean and standard deviation). At the end of the construction process, the chosen values of the temporary variables \mathbf{Z} are converted back into the original coordinate system \mathbf{X} .

IV. EXPERIMENTAL EVALUATION

An important collection of medical data sets for pattern classification is the well-known PROBEN1 data repository [18], which has been used various times in the past to evaluate and compare the performance of different methods for training NN classifiers. From the available data sets we chose *Cancer1*, *Diabetes1*, and *Heart1* for our experimentation.

Cancer1 concerns the diagnosis of breast cancer (possible outcomes: yes or no). The data set contains altogether 699 patterns. Each pattern has 9 input parameters (i.e., measurements). We used the first 525 of the patterns (i.e., about 75%) as training set (i.e., for optimizing the NN weights), and the remaining 174 as test set. *Diabetes1* concerns the diagnosis of diabetes (possible outcomes: yes or no). Each pattern has 8 input parameters. The data set contains altogether 768 patterns. We used the first 576 of them as training set and the remaining 192 as test set. Finally, *Heart1* concerns the diagnosis of a heart condition (possible outcomes: yes or no). Each pattern has 35 input parameters. The data set contains altogether 920 cases. We used the first 690 of them as training set and remaining 230 as test set.

Concerning the structure of the feed-forward NNs that we used, we took inspiration from the literature. More specifically we used the same network structures that were used in [1]. For an overview of these NN structures see Table I.

A. Algorithms for comparison

For comparison purposes we have re-implemented some algorithms traditionally used for training NNs, namely the back-propagation (BP) algorithm [19], and the Levenberg-Marquardt (LM) algorithm [15], [16]. Both algorithms (i.e., BP and LM) require gradient information, hence they require

the neuron transfer function $f(\cdot)$ to be differentiable. Consequently, these algorithms may not — in contrast to ACO* — be used in case the neuron transfer function is not differentiable or is unknown. In case of training NNs whose transfer function is differentiable, the drawback of general optimization algorithms such as ACO* is however that they do not exploit available additional information e.g., the gradient. In order to see how the additional gradient information influences the performance of ACO*, we have also implemented hybridized versions of ACO*, namely ACO*-BP and ACO*-LM. In these hybrids, each solution generated by the ACO* algorithm is improved by running a single improving iteration of either BP or LM before being evaluated.

Finally, we wanted to see how all the algorithms tested compare to a simple random search (RS) method. This is an algorithm that randomly generates a set of values for the weights and then evaluates these solutions. As we used a sigmoid function as neuron transfer function, it was sufficient to limit the range of weight values to values close to 0. Hence, we arbitrarily chose a range of [-5,5].

We have performed a limited scope parameter tuning for all algorithms using Birattari’s F-RACE method (see [4], [3]). The outcome is shown in Table II.

B. Results

In order to compare the performance of the algorithms, we applied each algorithm 50 times to each of the three test problems. As stopping condition we used the number of fitness function evaluations. Following the work of Alba and Chicano [1], we used 1000 function evaluations as the limit. Figures 3, 4, and 5 present respectively the results obtained for the cancer, diabetes, and heart test problems in the form of box-plots. Each figure presents the distributions of the actual classification error percentage (CEP) values obtained by the algorithms (over 50 independent runs).

Cancer1 (see Figure 3) appears to be the easiest data set among the three that we tackled. All algorithms obtained reasonably good results, including the RS method. However, the best performing algorithm is BP. From the fact that the results obtained by RS do not differ significantly from the results obtained by other—more complex algorithms, it may be concluded that the problem is relatively easy, and that there are a lot of reasonably good solutions scattered over the search space. None of the algorithms was able to classify all the test patterns correctly. This may be due to the limited size of the training set, i.e. there might have been not enough information in the training set to generalize perfectly.

Diabetes1 (see Figure 4) is a problem that is more difficult than *Cancer1*. All our algorithms clearly outperform RS. However, the overall performance of the algorithms in terms of the CEP value is not very good. The best performing is again BP. The less good overall performance of the algorithms may again indicate that the training set does not represent fully all the possible patterns.

TABLE II

SUMMARY OF THE FINAL PARAMETER VALUES THAT WE CHOSE FOR OUR ALGORITHMS. NOT INCLUDED IN THE TABLE ARE THE PARAMETERS COMMON TO ALL ACO* VERSIONS, NAMELY q AND m . FOR THESE PARAMETERS WE USED THE SETTINGS $q = 0.01$, AND $m = 2$ (THE NUMBER OF ANTS USED IN EACH ITERATION). NOTE THAT η IS THE STEP-SIZE PARAMETER OF BP, AND β IS THE ADAPTATION-STEP PARAMETER OF LM.

Algorithm	Cancer1				Diabetes1				Heart1			
	k	ρ	η	β	k	ρ	η	β	k	ρ	η	β
ACO*	148	0.95	-	-	136	0.8	-	-	230	0.6	-	-
ACO*-BP	148	0.98	0.3	-	136	0.7	0.1	-	230	0.98	0.4	-
ACO*-LM	148	0.9	-	10	136	0.1	-	10	230	0.1	-	10
BP	-	-	0.002	-	-	-	0.01	-	-	-	0.001	-
LM	-	-	-	50	-	-	-	5	-	-	-	1.5

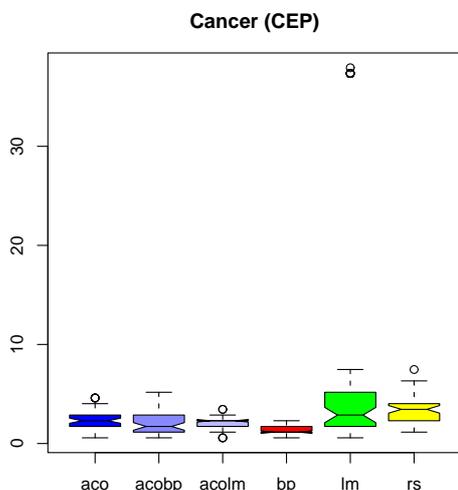


Fig. 3. Box-plots for Cancer1. The boxes are drawn between the first and the third quartile of the distribution, while the indentations in the box-plots (or notches) indicate the 95 % confidence interval.

The Heart1 problem (see Figure 5) is—with 230 weights—the largest problem that we tackled. It is also the one on which the performance of the algorithms differed mostly. All tested algorithms clearly outperform RS, but there are also significant differences among the more complex algorithms. BP, which was performing quite well on the other two test problems, did not do so well on Heart1. ACO* achieves results similar to BP. In turn, LM which was not performing so well on the first two problems, obtains quite good results. Very interesting is the performance of the hybridized versions of ACO*, namely ACO*-BP and ACO*-LM. The ACO*-BP hybrid clearly outperforms both ACO* and BP. ACO*-LM outperforms respectively ACO* and LM. Additionally, ACO*-LM performs best overall.

Finally, it is interesting to compare the performance of the ACO* based algorithms to some other general optimization algorithms. Alba and Chicano [1] have published the results of a genetic algorithm (GA) used for tackling exactly the same three problems as we did. They have tested not only a stand-alone GA, but also its hybridized versions: GA-BP

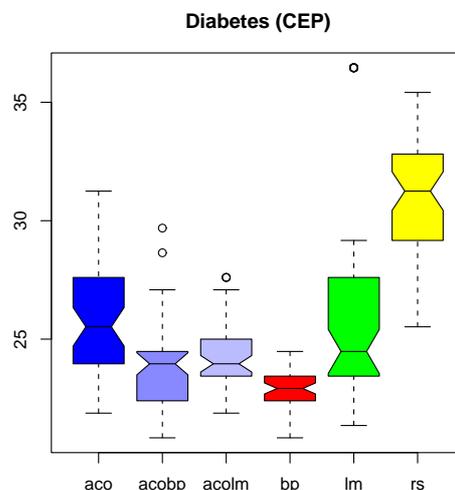


Fig. 4. Box-plots for Diabetes1. The boxes are drawn between the first and the third quartile of the distribution, while the indentations in the box-plots (or notches) indicate the 95 % confidence interval.

and GA-LM. Table III summarizes the results obtained by the ACO* and GA based algorithms. Clearly the stand-alone ACO* performs better than the stand-alone GA for all the test problems. ACO*-BP and ACO*-LM perform respectively better than GA-BP and GA-LM on both of the more difficult problems Diabetes1 and Heart1 and worse on Cancer1. For the Heart1 problem the mean performance of any ACO* based algorithm is significantly better than the best GA based algorithm (which was reported as the state-of-the-art for this problem in 2004).

V. CONCLUSION

We have presented an ant colony optimization algorithm (i.e., ACO*) for the training of feed-forward neural networks for pattern classification. The performance of the algorithm was evaluated on real-world test problems and compared to specialized algorithms for feed-forward neural network training (back propagation and Levenberg-Marquardt), and also to algorithms based on a genetic algorithm.

The performance of the stand-alone ACO* was comparable (or at least not much worse) than the performance of spe-

TABLE III

PAIR-WISE COMPARISON OF THE RESULTS OF THE ACO* BASED ALGORITHMS WITH RECENT RESULTS OBTAINED BY A SET OF GA BASED ALGORITHMS (SEE [1]). THE RESULTS CAN BE COMPARED DUE TO THE FACT THAT 1000 EVALUATIONS AS STOPPING CRITERION WERE USED FOR ALL THE ALGORITHMS. FOR EACH PROBLEM-ALGORITHM PAIR WE GIVE THE MEAN (OVER 50 INDEPENDENT RUNS), AND THE STANDARD DEVIATION (IN BRACKETS). THE BEST RESULT OF EACH COMPARISON IS INDICATED IN BOLD.

	GA	ACO*	GA-BP	ACO*-BP	GA-LM	ACO*-LM
Cancer1	16.76 (6.15)	2.39 (1.15)	1.43 (4.87)	2.14 (1.09)	0.02 (0.11)	2.08 (0.68)
Diabetes1	36.46 (0.00)	25.82 (2.59)	36.36 (0.00)	23.80 (1.73)	28.29 (1.15)	24.26 (1.40)
Heart1	41.50 (14.68)	21.59 (1.14)	54.30 (20.03)	18.29 (1.00)	22.66 (0.82)	16.53 (1.37)

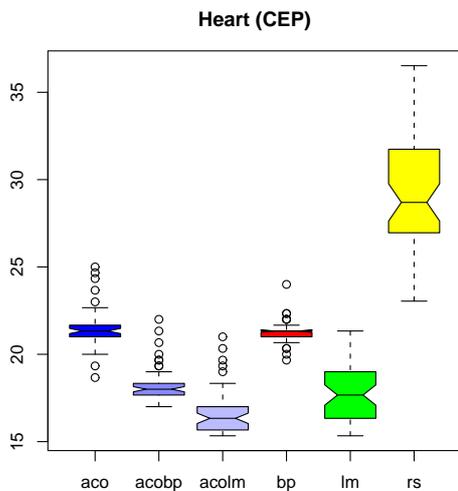


Fig. 5. Box-plots for Heart1. The boxes are drawn between the first and the third quartile of the distribution, while the indentations in the box-plots (or notches) indicate the 95 % confidence interval.

cialized algorithms for neural network training. This result is particularly interesting as ACO*—being a much more generic approach—allows also the training of networks in which the neuron transfer function is either not differentiable or unknown. The hybrid of ACO* and the Levenberg-Marquardt algorithm (i.e., ACO*-LM) was in some cases able to outperform the back propagation and the Levenberg-Marquardt algorithms. Finally, the results indicate that ACO* outperforms other general-purpose optimizers such as genetic algorithms.

ACKNOWLEDGMENT

This work was supported by the Spanish CICYT project TRACER (grant TIC-2002-04498-C05-03), and by the “Juan de la Cierva” program of the Spanish Ministry of Science and Technology of which Christian Blum is a post-doctoral research fellow. This work was also partially supported by the ANTS project, an Action de Recherche Concertée funded by the Scientific Research Directorate of the French Community of Belgium.

REFERENCES

- [1] E. Alba and J. F. Chicano, “Training neural networks with GA hybrid algorithms,” in *Proceedings of the Genetic and Evolutionary Computation Conference—GECCO 2004*, ser. Lecture Notes in Computer Science, K. D. et al., Ed., vol. 3102. Springer Verlag, Berlin, Germany, 2004, pp. 852–863.
- [2] B. Bilchev and I. C. Parmee, “The ant colony metaphor for searching continuous design spaces,” in *Proceedings of the AISB Workshop on Evolutionary Computation*, ser. Lecture Notes in Computer Science, vol. 993, 1995, pp. 25–39.
- [3] M. Birattari, “The problem of tuning metaheuristics as seen from a machine learning perspective,” Ph.D. dissertation, Université Libre de Bruxelles, Brussels, Belgium, 2004.
- [4] M. Birattari, T. Stützle, L. Paquete, and K. Varrentrapp, “A racing algorithm for configuring metaheuristics,” in *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2002)*, W. L. et al., Ed. Morgan Kaufmann Publishers, San Mateo, CA, 2002, pp. 11–18.
- [5] C. M. Bishop, *Neural networks for pattern recognition*. MIT Press, 2005.
- [6] E. Bonabeau, M. Dorigo, and G. Theraulaz, *Swarm Intelligence: From Natural to Artificial Systems*. Oxford University Press, New York, NY, 1999.
- [7] G. E. P. Box and M. E. Muller, “A note on the generation of random normal deviates,” *Annals of Mathematical Statistics*, vol. 29, no. 2, pp. 610–611, 1958.
- [8] J.-L. Deneubourg, S. Aron, S. Goss, and J.-M. Pasteels, “The self-organizing exploratory pattern of the argentine ant,” *Journal of Insect Behaviour*, vol. 3, pp. 159–168, 1990.
- [9] M. Dorigo, “Optimization, learning and natural algorithms (*in italian*),” Ph.D. dissertation, Dipartimento di Elettronica, Politecnico di Milano, Italy, 1992.
- [10] M. Dorigo, V. Maniezzo, and A. Coloni, “Ant System: Optimization by a colony of cooperating agents,” *IEEE Transactions on Systems, Man, and Cybernetics – Part B*, vol. 26, no. 1, pp. 29–41, 1996.
- [11] M. Dorigo and T. Stützle, *Ant Colony Optimization*. MIT Press, Cambridge, MA, 2004, to appear.
- [12] J. Dréo and P. Siarry, “A new ant colony algorithm using the heterarchical concept aimed at optimization of multimimima continuous functions,” in *Proceedings of ANTS 2002*, ser. Lecture Notes in Computer Science, M. Dorigo, G. Di Caro, and M. Sampels, Eds., vol. 2463. Springer Verlag, Berlin, Germany, 2002, pp. 216–221.
- [13] G. H. Golub and C. F. Loan, *Matrix Computations, 2nd ed.* The John Hopkins University Press, Baltimore, USA, 1989.
- [14] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer-Verlag, Berlin, Germany, 2001.
- [15] K. Levenberg, “A method for the solution of certain problems in least squares,” *Quarterly Applied Mathematics*, vol. 2, pp. 164–168, 1944.
- [16] D. Marquardt, “An algorithm for least-squares estimation of nonlinear parameters,” *SIAM Journal on Applied Mathematics*, vol. 11, pp. 431–441, 1963.
- [17] N. Monmarché, G. Venturini, and M. Slimane, “On how pachycondyla apicalis ants suggest a new search algorithm,” *Future Generation Computer Systems*, vol. 16, pp. 937–946, 2000.
- [18] L. Prechelt, “Proben1—a set of neural network benchmark problems and benchmarking rules,” Fakultät für Informatik, Universität Karlsruhe, Karlsruhe, Germany, Tech. Rep. 21, 1994.
- [19] D. Rummelhart, G. Hinton, and R. Williams, “Learning representations by backpropagation errors,” *Nature*, vol. 323, pp. 533–536, 1986.
- [20] K. Socha, “Extended ACO for continuous and mixed-variable optimization,” in *Proceedings of ANTS 2004 – Fourth International Workshop on Ant Algorithms and Swarm Intelligence*, ser. Lecture Notes in Computer Science, M. Dorigo, M. Birattari, C. Blum, L. M. Gambardella, F. Mondada, and T. Stützle, Eds. Springer Verlag, Berlin, Germany, 2004, to appear.