

Brief Papers

Training Feedforward Networks with the Marquardt Algorithm

Martin T. Hagan and Mohammad B. Menhaj

Abstract— The Marquardt algorithm for nonlinear least squares is presented and is incorporated into the backpropagation algorithm for training feedforward neural networks. The algorithm is tested on several function approximation problems, and is compared with a conjugate gradient algorithm and a variable learning rate algorithm. It is found that the Marquardt algorithm is much more efficient than either of the other techniques when the network contains no more than a few hundred weights.

I. INTRODUCTION

SINCE the backpropagation learning algorithm [1] was first popularized, there has been considerable research on methods to accelerate the convergence of the algorithm. This research falls roughly into two categories. The first category involves the development of ad hoc techniques (e.g., [2]–[5]). These techniques include such ideas as varying the learning rate, using momentum and rescaling variables. Another category of research has focused on standard numerical optimization techniques (e.g., [6]–[9]).

The most popular approaches from the second category have used conjugate gradient or quasi-Newton (secant) methods. The quasi-Newton methods are considered to be more efficient, but their storage and computational requirements go up as the square of the size of the network. There have been some limited memory quasi-Newton (one step secant) algorithms that speed up convergence while limiting memory requirements [8,10]. If exact line searches are used, the one step secant methods produce conjugate directions.

Another area of numerical optimization that has been applied to neural networks is nonlinear least squares [11]–[13]. The more general optimization methods were designed to work effectively on all sufficiently smooth objective functions. However, when the form of the objective function is known it is often possible to design more efficient algorithms. One particular form of objective function that is of interest for neural networks is a sum of squares of other nonlinear functions. The minimization of objective functions of this type is called nonlinear least squares.

Most of the applications of nonlinear least squares to neural networks have concentrated on sequential implementations, where the weights are updated after each presentation of an input/output pair. This technique is useful when on-line adaptation is needed, but it requires that several approximations be made to the standard algorithms. The standard algorithms are performed in batch mode, where the weights are only updated after a complete sweep through the training set.

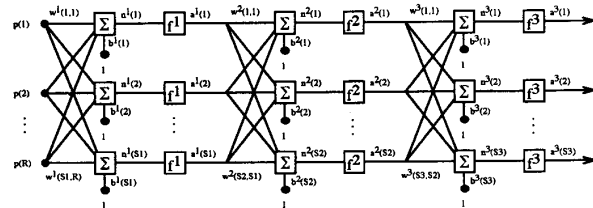


Fig. 1. Three-layer feedforward network.

This paper presents the application of a nonlinear least squares algorithm to the batch training of multi-layer perceptrons. For very large networks the memory requirements of the algorithm make it impractical for most current machines (as is the case for the quasi-Newton methods). However, for networks with a few hundred weights the algorithm is very efficient when compared with conjugate gradient techniques. Section II briefly presents the basic backpropagation algorithm. The main purpose of this section is to introduce notation and concepts which are needed to describe the Marquardt algorithm. The Marquardt algorithm is then presented in Section III. In Section IV the Marquardt algorithm is compared with the conjugate gradient algorithm and with a variable learning rate variation of backpropagation. Section V contains a summary and conclusions.

II. BACKPROPAGATION ALGORITHM

Consider a multilayer feedforward network, such as the three-layer network of Fig. 1.

The net input to unit i in layer $k + 1$ is

$$n^{k+1}(i) = \sum_{j=1}^{S_k} w^{k+1}(i, j)a^k(j) + b^{k+1}(i). \quad (1)$$

The output of unit i will be

$$a^{k+1}(i) = f^{k+1}(n^{k+1}(i)). \quad (2)$$

For an M layer network the system equations in matrix form are given by

$$\underline{a}^0 = \underline{p} \quad (3)$$

$$\underline{a}^{k+1} = \underline{f}^{k+1}(W^{k+1}\underline{a}^k + \underline{b}^{k+1}), \quad k = 0, 1, \dots, M - 1. \quad (4)$$

The task of the network is to learn associations between a specified set of input–output pairs $\{(p_1, t_1), (p_2, t_2), \dots, (p_Q, t_Q)\}$.

The performance index for the network is

$$V = \frac{1}{2} \sum_{q=1}^Q (t_q - a_q^M)^T (t_q - a_q^M) = \frac{1}{2} \sum_{q=1}^Q e_q^T e_q \quad (5)$$

where a_q^M is the output of the network when the q th input, p_q , is presented, and $e_q = t_q - a_q^M$ is the error for the q th input. For the standard backpropagation algorithm we use an approximate steepest descent rule. The performance index is approximated by

$$\hat{V} = \frac{1}{2} e_q^T e_q \quad (6)$$

where the total sum of squares is replaced by the squared errors for a single input/output pair. The approximate steepest (gradient) descent algorithm is then

$$\Delta w^k(i, j) = -\alpha \frac{\partial \hat{V}}{\partial w^k(i, j)} \quad (7)$$

$$\Delta b^k(i) = -\alpha \frac{\partial \hat{V}}{\partial b^k(i)} \quad (8)$$

where α is the learning rate. Define

$$\delta^k(i) \equiv \frac{\partial \hat{V}}{\partial n^k(i)} \quad (9)$$

as the sensitivity of the performance index to changes in the net input of unit i in layer k . Now it can be shown, using (1), (6), and (9), that

$$\frac{\partial \hat{V}}{\partial w^k(i, j)} = \frac{\partial \hat{V}}{\partial n^k(i)} \frac{\partial n^k(i)}{\partial w^k(i, j)} = \delta^k(i) a^{k-1}(j) \quad (10)$$

$$\frac{\partial \hat{V}}{\partial b^k(i)} = \frac{\partial \hat{V}}{\partial n^k(i)} \frac{\partial n^k(i)}{\partial b^k(i)} = \delta^k(i). \quad (11)$$

It can also be shown that the sensitivities satisfy the following recurrence relation

$$\underline{\delta}^k = \hat{F}^k(\underline{n}^k) W^{k+1 \top} \underline{\delta}^{k+1} \quad (12)$$

where

$$\hat{F}^k(\underline{n}^k) = \begin{bmatrix} f^k(n^k(1)) & 0 & \dots & 0 \\ 0 & f^k(n^k(2)) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f^k(n^k(S_k)) \end{bmatrix} \quad (13)$$

and

$$f^k(n) = \frac{df^k(n)}{dn}. \quad (14)$$

This recurrence relation is initialized at the final layer

$$\underline{\delta}^M = -\hat{F}^M(\underline{n}^M)(t_q - a_q). \quad (15)$$

The overall learning algorithm now proceeds as follows; first, propagate the input forward using (3)–(4); next, propagate the sensitivities back using (15) and (12); and finally, update the weights and offsets using (7), (8), (10), and (11).

III. MARQUARDT-LEVENBERG MODIFICATION

While backpropagation is a steepest descent algorithm, the Marquardt-Levenberg algorithm [14] is an approximation to Newton's method. Suppose that we have a function $V(\underline{x})$ which we want to minimize with respect to the parameter vector \underline{x} , then Newton's method would be

$$\Delta \underline{x} = -[\nabla^2 V(\underline{x})]^{-1} \nabla V(\underline{x}) \quad (16)$$

where $\nabla^2 V(\underline{x})$ is the Hessian matrix and $\nabla V(\underline{x})$ is the gradient. If we assume that $V(\underline{x})$ is a sum of squares function

$$V(\underline{x}) = \sum_{i=1}^N e_i^2(\underline{x}) \quad (17)$$

then it can be shown that

$$\nabla V(\underline{x}) = J^T(\underline{x}) \underline{e}(\underline{x}) \quad (18)$$

$$\nabla^2 V(\underline{x}) = J^T(\underline{x}) J(\underline{x}) + S(\underline{x}) \quad (19)$$

where $J(\underline{x})$ is the Jacobian matrix

$$J(\underline{x}) = \begin{bmatrix} \frac{\partial e_1(\underline{x})}{\partial x_1} & \frac{\partial e_1(\underline{x})}{\partial x_2} & \dots & \frac{\partial e_1(\underline{x})}{\partial x_n} \\ \frac{\partial e_2(\underline{x})}{\partial x_1} & \frac{\partial e_2(\underline{x})}{\partial x_2} & \dots & \frac{\partial e_2(\underline{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e_N(\underline{x})}{\partial x_1} & \frac{\partial e_N(\underline{x})}{\partial x_2} & \dots & \frac{\partial e_N(\underline{x})}{\partial x_n} \end{bmatrix} \quad (20)$$

and

$$S(\underline{x}) = \sum_{i=1}^N e_i(\underline{x}) \nabla^2 e_i(\underline{x}). \quad (21)$$

For the Gauss-Newton method it is assumed that $S(\underline{x}) \approx 0$, and the update (16) becomes

$$\Delta \underline{x} = [J^T(\underline{x}) J(\underline{x})]^{-1} J^T(\underline{x}) \underline{e}(\underline{x}). \quad (22)$$

The Marquardt-Levenberg modification to the Gauss-Newton method is

$$\Delta \underline{x} = [J^T(\underline{x}) J(\underline{x}) + \mu I]^{-1} J^T(\underline{x}) \underline{e}(\underline{x}). \quad (23)$$

The parameter μ is multiplied by some factor (β) whenever a step would result in an increased $V(\underline{x})$. When a step reduces $V(\underline{x})$, μ is divided by β . (In Section IV we used $\mu = 0.01$ as a starting point, with $\beta = 10$.) Notice that when μ is large the algorithm becomes steepest descent (with step $1/\mu$), while for small μ the algorithm becomes Gauss-Newton. The Marquardt-Levenberg algorithm can be considered a trust-region modification to Gauss-Newton [8].

The key step in this algorithm is the computation of the Jacobian matrix. For the neural network mapping problem the terms in the Jacobian matrix can be computed by a simple modification to the backpropagation algorithm. The performance index for the mapping problem is given by (5). It is easy to see that this is equivalent in form to (17), where $\underline{x} = [w^1(1, 1)w^1(1, 2) \dots w^1(S_1, R)b^1(1) \dots b^1(S_1)w^2(1, 1) \dots b^M(SM)]^T$, and $N = Q \times SM$. Standard backpropagation calculates terms like

$$\frac{\partial \hat{V}}{\partial w^k(i, j)} = \frac{\partial \sum_{m=1}^{SM} e_q^2(m)}{\partial w^k(i, j)}. \quad (24)$$

For the elements of the Jacobian matrix that are needed for the Marquardt algorithm we need to calculate terms like

$$\frac{\partial e_q(m)}{\partial w^k(i,j)} \quad (25)$$

These terms can be calculated using the standard backpropagation algorithm with one modification at the final layer

$$\Delta^M = -\hat{F}^M(\underline{n}^M). \quad (26)$$

Note that each column of the matrix in (26) is a sensitivity vector which must be backpropagated through the network to produce one row of the Jacobian.

A. Summary

The Marquardt modification to the backpropagation algorithm thus proceeds as follows:

- 1) Present all inputs to the network and compute the corresponding network outputs (using (3) and (4)), and errors ($e_q = t_q - a_q^M$). Compute the sum of squares of errors over all inputs ($V(\underline{x})$).
- 2) Compute the Jacobian matrix (using (26), (12), (10), (11), and (20)).
- 3) Solve (23) to obtain $\Delta \underline{x}$. (For the results shown in the next section Cholesky factorization was used to solve this equation.)
- 4) Recompute the sum of squares of errors using $\underline{x} + \Delta \underline{x}$. If this new sum of squares is smaller than that computed in step 1, then reduce μ by β , let $\underline{x} = \underline{x} + \Delta \underline{x}$, and go back to step 1. If the sum of squares is not reduced, then increase μ by β and go back to step 3.
- 5) The algorithm is assumed to have converged when the norm of the gradient ((18)) is less than some predetermined value, or when the sum of squares has been reduced to some error goal.

IV. RESULTS

The Marquardt backpropagation algorithm (MBP), as described in the previous section, was tested on five function approximation problems. To provide a basis for comparison, two other modifications to backpropagation were also applied to the same problems: backpropagation with variable learning rate (VLBP) [2], and conjugate gradient backpropagation (CGBP). For the purposes of this study we used the Fletcher-Reeves version of the conjugate gradient algorithm [15, pp. 73–84], with an exact line search. The line search consisted of two parts: interval location, using function comparison [15, pp. 41–42] and a golden section search [15, pp. 31–32]. It should be noted that there are a number of decisions to be made in the implementation of the conjugate gradient algorithm, including the precision and the type of line search and the choice of the number of steps before the search direction is reinitialized to the gradient direction (typically chosen to be equal to the number of parameters, see [15]). We took some care in making these decisions, but there is no guarantee that our implementation is optimal. However, we found that the basic trends described in the examples to follow were

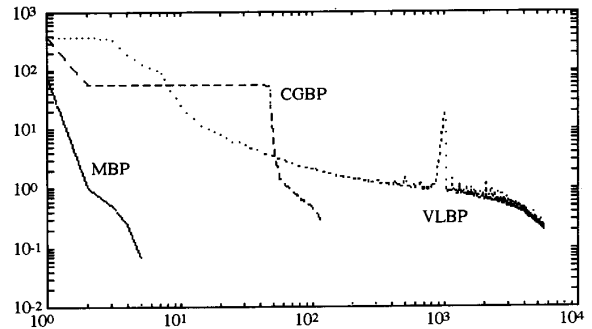


Fig. 2. Network convergence for sine wave (sum of squares vs. epoch).

not sensitive to modifications in the implementation of the conjugate gradient algorithm.

A. Problem #1: Sine Wave

For the first test problem a 1–15–1 network, with a hidden layer of sigmoid nonlinearities and a linear output layer, was trained to approximate the sinusoidal function

$$y = 1/2 + 1/4 \sin(3\pi x)$$

using MBP, CGBP and VLBP. Fig. 2 displays the training curves for the three methods. The training set consisted of 40 input/output pairs, where the input values were scattered in the interval $[-1,1]$; and the network was trained until the sum of squares of the errors was less than the error goal of 0.02. The curves shown in Fig. 2 are an average over 5 different initial weights. The initial weights are random, but are normalized using the method of Nguyen and Widrow [16].

A comment is appropriate at this point on the shape of the learning curve for CGBP. Note that there is a plateau followed by a steep curve. The plateau ends at the point when the search direction is re-initialized to the gradient direction (when the number of steps is equal to the number of parameters), and the plateau only occurs once. When we reset the search direction more often the plateau was shortened, but the steepness of the subsequent slope was reduced, and the overall convergence rate was not improved. We found that the plateau could be eliminated by using the gradient direction for the first few steps and then using the conjugate gradient algorithm, but this made only a small difference in overall convergence rate.

Fig. 2 provides only limited information, since the three algorithms do not have the same number of floating point operations for each iteration. The first line of Table 1 summarizes the results, showing the number of floating point operations required for convergence. Notice that CGBP takes more than nine times as many flops as MBP, and VLBP takes almost 45 times as many flops.

B. Problem #2: Square Wave

Fig. 3 illustrates the second test problem, in which the same 1–15–1 network is trained to approximate a square wave. Fig. 4 displays the average learning curves (5 different initial weights), and line 2 of Table I summarizes the average results.

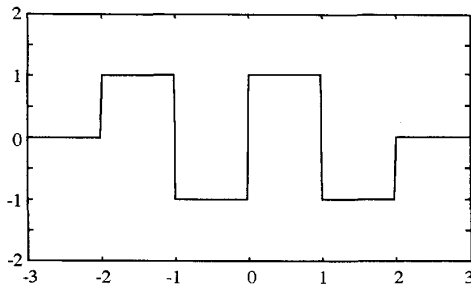


Fig. 3. Function approximation, problem #2.

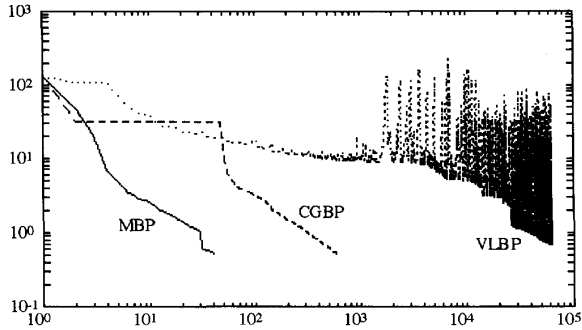


Fig. 4. Network convergence for square wave (sum of squares vs. epoch).

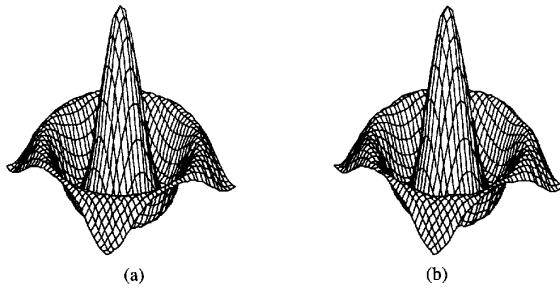


Fig. 5. Function approximation, problem #3, and network response; (a) sinc function, (b) network response.

Note that CGBP takes more than four times as many flops as MBP, and VLBP takes more than 65 times as many flops.

C. Problem #3: 2-D Sinc Function

Fig. 5(a) illustrates the third test problem. In this case a 2-15-1 network is trained to approximate a two-dimensional sinc function. Fig. 6 displays the average learning curves (3 different initial weights), and line 3 of Table I summarizes the average results (error goal of 0.5 with 289 input/output sets). The CGBP algorithm takes more than 7 times as many flops as MBP, and VLBP takes more than 27 times as many flops. These differences in convergence time became more pronounced as the error goal was reduced, but the time required for convergence of VLBP made multiple runs impractical. Fig. 5(b) illustrates the network response, after training with MBP to an error goal of 0.02.

TABLE II
NUMBER OF FLOPS REQUIRED FOR CONVERGENCE

	VLBP	CGBP	MBP
Sine Wave	8.42×10^7	1.75×10^7	1.89×10^6
Square Wave	2.28×10^9	1.49×10^8	3.48×10^7
2-D Sinc	2.94×10^9	7.67×10^8	1.07×10^8
4-D Test	—	7.71×10^9	1.97×10^9

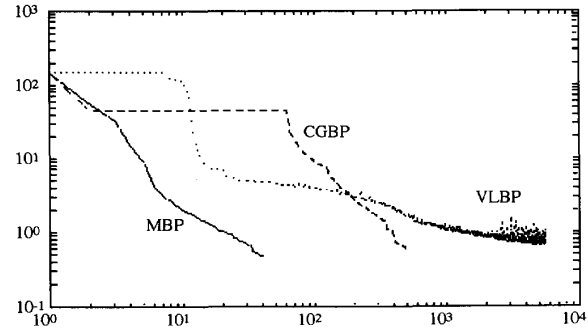


Fig. 6. Network convergence for 2-D sinc function (sum of squares vs. epoch).

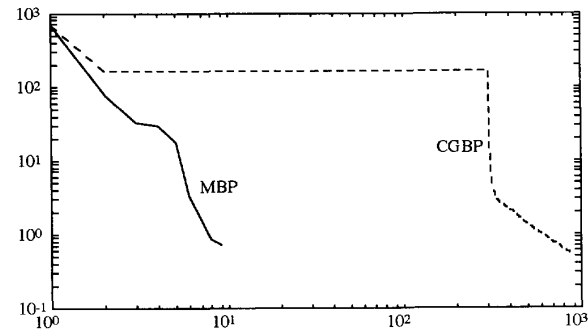


Fig. 7. Network convergence for 4-D test function (sum of squares vs. epoch).

D. Problem #4: 4-D Function

The fourth test problem is a four input-single output function

$$y = \sin(2\pi x_1) x_2^2 x_3^3 x_4^4 e^{-(x_1+x_2+x_3+x_4)}. \quad (27)$$

For this example a 4-50-1 network (301 parameters) is trained to approximate (27), where the input values were scattered in the interval [-1,1]; and the network was trained until the sum of squares of the errors (over 400 input/output sets) was less than the error goal of 0.5. Fig. 7 displays the average learning curves (3 different initial weights), and line 4 of Table I summarizes the average results. The CGBP algorithm takes approximately four times as many flops as MBP (VLBP was not applied to this problem because of excessive training time)

E. Example Set Scaling

In order to investigate the effect of sample size on the efficiency of the algorithm, MBP and CGBP were used to

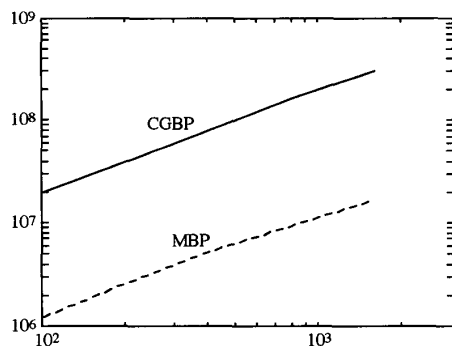


Fig. 8. FLOPS required for convergence versus sample size.

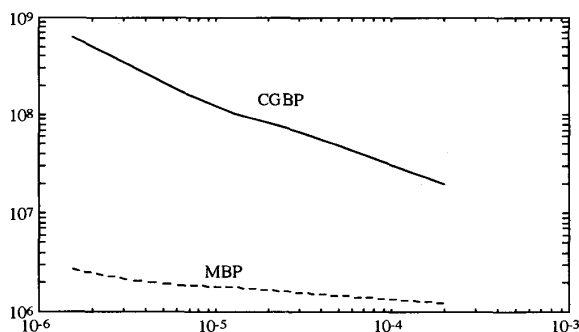


Fig. 9. FLOPS required for convergence versus error goal.

train a 1-10-1 network to approximate the function

$$y = \sin(\pi/2x) \quad (28)$$

over the interval $-1 < x < 1$. The size of the training set was varied from 100 points to 1600 points, and the network was trained until the mean square error was less than 2×10^{-4} . Fig. 8 displays the number of flops required for convergence, as a function of the sample size. The curves represent an average over 10 different initial conditions. From this figure we can see that the effect is linear, both for MBP and CGBP. MBP is approximately 16 times faster than CGBP for each sample size.

F. Accuracy Requirements

We noted that the difference between the performances of MBP and CGBP became more pronounced as higher precision approximations were required. This effect is illustrated in Fig. 9. In this example a 1-10-1 network is trained to approximate the sine wave of (28). The sample size is held constant at 100 points, but the mean square error goal is halved in steps from 2×10^{-4} to 1.6×10^{-6} . Fig. 9 displays the number of flops required for convergence, as a function of the error goal, for both MBP and CGBP. The curves represent an average over 10 different initial conditions. With an error goal of 2×10^{-4} MBP is 16 times faster than CGBP. This ratio increases as the error goal is reduced; when the error goal is 1.6×10^{-6} , MBP is 136 times faster than CGBP.

V. CONCLUSION

Many numerical optimization techniques have been successfully used to speed up convergence of the backpropagation learning algorithm. This paper presented a standard nonlinear least squares optimization algorithm, and showed how to incorporate it into the backpropagation algorithm. The Marquardt algorithm was tested on several function approximation problems, and it was compared with the conjugate gradient algorithm and with variable learning rate backpropagation. The results indicate that the Marquardt algorithm is very efficient when training networks which have up to a few hundred weights. Although the computational requirements are much higher for each iteration of the Marquardt algorithm, this is more than made up for by the increased efficiency. This is especially true when high precision is required.

The authors also found that in many cases the Marquardt algorithm converged when the conjugate gradient and variable learning rate algorithms failed to converge. For example, in problem #2 (Fig. 3) if we used five neurons in the hidden layer the CGBP and VLBP algorithms almost never converged to an optimal solution. The MBP algorithm converged to an optimal solution in 50% of the tests, and in less time than was required for the network with 15 hidden neurons.

REFERENCES

- [1] D. E. Rumelhart, G. E. Hinton and R. J. Williams, "Learning Representations by Back-Propagating Errors," *Nature*, vol. 323, pp. 533-536, 1986.
- [2] T. P. Vogl, J. K. Mangis, A. K. Zigler, W. T. Zink and D. L. Alkon, "Accelerating the convergence of the backpropagation method," *Bio. Cybern.*, vol. 59, pp. 256-264, Sept. 1988.
- [3] R. A. Jacobs, "Increased Rates of Convergence Through Learning Rate Adaptation," *Neural Networks*, vol. 1, no. 4, pp. 295-308, 1988.
- [4] T. Tollenaere, "SuperSAB: Fast adaptive back propagation with good scaling properties," *Neural Networks*, vol. 3, no. 5, pp. 561-573, 1990.
- [5] A. K. Rigler, J. M. Irvine, and T. P. Vogl, "Rescaling of variables in back propagation learning," *Neural Networks*, vol. 3, no. 5, pp. 561-573, 1990.
- [6] D. F. Shanno, "Recent advances in numerical techniques for large-scale optimization," in *Neural Networks for Control*, Miller, Sutton and Werbos, Eds. Cambridge MA: MIT Press, 1990.
- [7] E. Barnard, "Optimization for training neural nets," *IEEE Trans. Neural Net.*, vol. 3, no. 2, pp. 232-240, 1992.
- [8] R. Battiti, "First- and second order methods for learning: Between steepest descent and Newton's method," *Neural Computation*, vol. 4, no. 2, pp. 141-166, 1992.
- [9] C. Charalambous, "Conjugate gradient algorithm for efficient training of artificial neural networks," *IEE Proc.*, vol. 139, no. 3, pp. 301-310, 1992.
- [10] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Math. Prog.*, vol. 45, pp. 503-528, 1989.
- [11] S. Kollias and D. Anastassiou, "An adaptive least squares algorithm for the efficient training of artificial neural networks," *IEEE Trans. Circ. Syst.*, vol. 36, no. 8, pp. 1092-1101, 1989.
- [12] S. Singhal and L. Wu, "Training multilayer perceptrons with the Extended Kalman Algorithm," in *Advances in Neural Information Processing Systems I*, D. S. Touretzky, Ed. San Mateo, CA: Morgan Kaufman, pp. 133-140, 1989.
- [13] G. V. Puskorius and L. A. Feldkamp, "Decoupled Extended Kalman Filter Training of feedforward layered networks," in *Proc. IJCNN*, vol. 1, pp. 771-777, July 1991.
- [14] D. Marquardt, "An algorithm for least squares estimation of non-linear parameters," *J. Soc. Ind. Appl. Math.*, pp. 431-441, 1963.
- [15] L. E. Scales, *Introduction to Nonlinear Optimization*. New York: Springer-Verlag, 1985.
- [16] D. Nguyen and B. Widrow, "Improving the learning speed of 2-layer neural networks by choosing initial values of the adaptive weights," in *Proc. IJCNN*, vol. 3, pp. 21-26, July 1990.