

# Efficient Numerical Inversion Using Multilayer Feedforward Neural Networks

**George Lendaris**

Professor for Systems Science and Electrical Engineering  
Portland State University  
Portland, Oregon  
lendaris@syc.pdx.edu

**Karl Mathia**

Senior Research Engineer  
Accurate Automation Corporation  
Chattanooga, Tennessee  
kmathia@accurate-automation.com

## Abstract

*An efficient second-order optimization algorithm for the numerical inversion of nonlinear functions using feedforward neural networks is presented. After a function has been successfully learned by a neural network, the (nonlinear) equation is solved using the neural network representation, in effect, performing a numerical inversion process. As is well known, second-order information can dramatically improve the convergence of numerical methods. It is demonstrated, that the algorithm developed herein, based on a neural network, also dramatically reduces the amount of computations required, hence yielding a more efficient numerical inversion process. Certain constraints must also be considered. We discuss these issues, present an example and briefly characterize an applicable class of problems.*

## 1 Introduction

An efficient second-order optimization algorithm for the numerical inversion of nonlinear functions using feedforward neural networks is presented. After a function has been successfully learned by a neural network, the (nonlinear) equation is solved using the neural network representation, in effect, performing a numerical inversion process. For example, this is useful when a function is given only in terms of sample points. A common scenario for solving nonlinear and high-dimensional functions is to define and solve an optimization problem. We follow this strategy and define a scalar, multi-variate error function  $E(\mathbf{x})$  and search for the optimal solution vector  $\mathbf{x}^*$ , such that  $E(\mathbf{x}^*) = \min\{E(\mathbf{x})\}$ . Many numerical optimization schemes are available for this purpose. A technique is called of 'order  $n$ ' if it uses a  $n$ -th order Taylor series expansion to approximate  $E$ , and then iteratively minimizes and updates the approximating function instead of the original. Incorporating higher order information, i.e. up to the  $n$ -th derivative of  $E$ , provides a better approximation and therefore can dramatically improve convergence. Unfortunately, higher order terms also can make the numerical process unreasonably complex. Many optimization schemes therefore approximate not only the original function, but also its derivatives. We design and evaluate an efficient algorithm for a class of feedforward neural networks. Some important limitations of the algorithm are also discussed.

In the following sections we define the architecture of a feedforward multilayer perceptron (MLP) which represents the family of neural networks under consideration. Second-order techniques of decreasing computational complexity are reviewed, 'decreasing' in the sense that - depending on the algorithm - derivatives and matrix inversions are computed or approximated. We apply the algorithms to the networks previously defined, and demonstrate the techniques (and its limitations) with an example. We conclude with a discussion and characterization of an applicable class of problems.

## 2 Feedforward Neural Networks

The class of neural networks used as an example in this paper is the common MLP, although the algorithm is applicable to similar feedforward architectures with only minor modifications.

### 2.1 Universal Approximator Network

The feedforward neural network defined in the universal approximator theorem [4], the MLP, is used as an example network for design and evaluation of the algorithm. The network has  $n$  input nodes, one

hidden layer with  $q$  sigmoidal processing elements (PEs), and an output layer with  $n$  linear PEs. The network represents a nonlinear mapping  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  with input vector  $\mathbf{x}$  and output vector  $\mathbf{y}$ , defined by

$$\mathbf{f}: \mathcal{R}^n \rightarrow \mathcal{R}^n, \quad \mathbf{f}(\mathbf{x}) = \mathbf{W} \cdot \sigma(\mathbf{V}\mathbf{x} + \mathbf{b}) = \mathbf{W} \cdot \sigma(\mathbf{a}). \quad (1)$$

After training, the weight matrices of the hidden and output layers,  $\mathbf{V} \in \mathcal{R}^{q \times n}$  and  $\mathbf{W} \in \mathcal{R}^{n \times q}$ , are fixed. The usual bias input  $\mathbf{b}$  is applied to the hidden layer only. We define  $\mathbf{a} = \mathbf{V}\mathbf{x} + \mathbf{b}$ . The (column) vector-valued function  $\sigma(\mathbf{a}) = [\sigma(a_1), \dots, \sigma(a_q)]^T$  represents the sigmoid transfer functions of PEs in the hidden layer, where  $\sigma(a_i) = 1/(1 + \exp(-a_i))$  is the output of the  $i$ -th PE.

## 2.2 Jacobian and Hessian

The Jacobian (matrix of first derivatives) and the Hessian (matrix of second derivatives) of the MLP in Equation 1 is needed for the derivation of the algorithms in Section 4. The Jacobian is

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{W} \cdot \frac{\partial \sigma(\mathbf{a})}{\partial \mathbf{a}} \cdot \mathbf{V} = \mathbf{W} \cdot \Sigma'(\mathbf{a}) \cdot \mathbf{V}, \quad (2)$$

where  $\Sigma'(\mathbf{a})$  is a diagonal matrix with elements  $\sigma'(a_i) = \sigma(a_i) \cdot (1 - \sigma(a_i))$ , the well known form of a sigmoid's derivative. The Hessian's element in row  $i$  and column  $j$  is

$$\frac{\partial^2}{\partial x_i \partial x_j} \mathbf{f}(\mathbf{x}) = \mathbf{W} \cdot (\Sigma''(\mathbf{a}) \cdot \mathbf{v}^i) \cdot \mathbf{v}^j, \quad (3)$$

where  $\mathbf{v}^i$  and  $\mathbf{v}^j$  are the  $i$ -th and  $j$ -th column vectors in  $\mathbf{V}$ , respectively.  $\Sigma''(\mathbf{a})$  is a diagonal matrix with elements  $\sigma''(a_i) = \sigma(a_i) - 3\sigma^2(a_i) + 2\sigma^3(a_i)$ . It is clear that the computation of the Hessian (second-order information!) even for this relatively simple network architecture is considerably complex [5].

## 3 First-Order Optimization

A popular first-order technique in the context of artificial neural networks is error backpropagation [9], a training algorithm for optimizing the connection weights of a feedforward ANN for a given problem. 'Backprop' can be viewed as an unconstrained nonlinear optimization scheme which combines the classical gradient descent and backsubstitution methods. It is based on the assumption that a first-order approximation of  $E$  about some point  $\mathbf{x}$  is sufficiently accurate, thus  $E(\mathbf{x} + \Delta\mathbf{x}) = E(\mathbf{x}) + E'(\mathbf{x})\Delta\mathbf{x}$ , where  $E' = \partial E / \partial \mathbf{x}$  is the error gradient. Backprop iteratively minimizes and updates the approximating function. A potential drawback of first-order methods is poor convergence, which mainly corresponds to a poor first-order approximation. The backprop weight update rule therefore often includes the 'momentum term',

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta\mathbf{x}_k + \Delta\mathbf{x}_{k-1} = \mathbf{x}_k + \eta[E'(\mathbf{x})]^T + \mu(\mathbf{x}_k - \mathbf{x}_{k-1}), \quad (4)$$

where  $\mathbf{x}$  denotes a vector of connection weights and the scalar  $\eta$  the learning rate. The momentum term  $\Delta\mathbf{x}_{k-1}$ , with the momentum constant  $\mu$ , has been suggested as a (not always successful) *ad hoc* attempt to include *second-order information* in order to increase convergence speed. Also note that training algorithms like backprop search the *weight space* of a network, whereas here we are concerned with searching the *input space* (see also [6]).

## 4 Second-Order Optimization

In this section we review second-order optimization techniques of decreasing computational complexity, decreasing in the sense that - depending on the algorithm - the Hessian of a scalar, multi-variate error function is 1) computed *and* inverted (Newton's method), 2) *only* computed and its inversion avoided (con-

jugate gradient algorithm), 3) *both* computation and inversion of the Hessian are avoided (scaled conjugate gradient algorithm). We focus on 3), but also perform the computations in 1) and 2) for comparison purposes (see Section 5.2).

Second-order optimization techniques assume that a second-order Taylor series expansion is a sufficiently accurate approximation of the error function  $E$  about a point  $\mathbf{x}$ ,

$$E(\mathbf{x} + \Delta\mathbf{x}) = E(\mathbf{x}) + \Delta E(\mathbf{x}) = E(\mathbf{x}) + E'(\mathbf{x}) \cdot \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T \cdot E''(\mathbf{x}) \cdot \Delta\mathbf{x}, \quad (5)$$

where  $E''$  is the Hessian of  $E$ . This assumption is appropriate near minima, where higher order terms vanish. Of course we assume that  $E'$  and  $E''$  exist and are continuous. The approximating function is minimized by solving

$$\frac{\partial}{\partial \Delta\mathbf{x}} \{ \Delta E(\mathbf{x}) \} = E'(\mathbf{x}) + \Delta\mathbf{x}^T \cdot E''(\mathbf{x}) = \mathbf{0}^T, \quad (6)$$

with the solution

$$\Delta\mathbf{x}^* = -[E''(\mathbf{x})]^{-1} \cdot [E'(\mathbf{x})]^T. \quad (7)$$

The Hessian  $E''$  must be *positive definite* [3]. This is a stringent constraint as we will demonstrate below. Also note that the Hessian of continuous functions is symmetric, and that the gradient  $E'$  is a row vector [2]. For the remainder of this paper we write  $\mathbf{g} = E'$  and  $\mathbf{H} = E''$ .

#### 4.1 Newton's Method

The key calculation for all second-order optimization techniques based on Newton's methods is Equation 7. With  $\Delta\mathbf{x} = \mathbf{x}_{k+1} - \mathbf{x}_k$  and our new notation, Newton's method is the iteration process

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_k^{-1} \cdot \mathbf{g}_k^T. \quad (8)$$

with  $\mathbf{H}_k = \mathbf{H}(\mathbf{x}_k)$ ,  $\mathbf{g}_k = \mathbf{g}(\mathbf{x}_k)$ . One iteration step minimizes the approximating function in Equation 5. Many iterations may be needed until the minimum of the approximating function is sufficiently close to the minimum of the original function. A drawback of Newton's method and its variants is the computation, storage, and inversion of the Hessian at each iteration step, which can be computationally expensive, in particular for high-dimensional systems. Quasi-Newton methods avoid these computations, for example the conjugate gradient algorithm.

#### 4.2 Conjugate Gradient Method

The conjugate gradient (CG) method iteratively approximates the inverse Hessian matrix, while limiting the number of iteration steps to  $n$  by searching only along basis vectors of the  $n$ -dimensional search space. The vector from starting point  $\mathbf{x}_1$  to the solution  $\mathbf{x}^*$  is therefore some linear combination of  $n$  basis vectors  $\mathbf{p}_i$  (row vectors) of length  $\alpha_i$ ,

$$\mathbf{x}^* - \mathbf{x}_1 = \sum_{i=1}^n \alpha_i \cdot \mathbf{p}_i. \quad (9)$$

The question is: How can we obtain an appropriate set of basis vectors and the associated step sizes? Fortunately, if the gradient and the Hessian are known at the starting point, these entities can be computed while the algorithm progresses, based on the orthogonality of each new gradient  $\mathbf{g}_{k+1}$  to the previous search direction  $\mathbf{p}_k$  [3]. Which particular basis is being constructed (one out of an infinite number of possible bases) depends on the starting point  $\mathbf{x}_1$ . The standard CG algorithm is shown in Figure 1.

Choose  $\mathbf{x}_1$  (column vector). Compute  $\mathbf{H} = \mathbf{E}''(\mathbf{x}_1)$ ,  $\mathbf{g}_1 = -\mathbf{E}'(\mathbf{x}_1)$ . Set  $\mathbf{p}_1 = \mathbf{g}_1^T$  (column vector).  
 LOOP:  $i = 1..\text{dim}(\mathbf{x})$

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{p}_i, \quad \alpha_i = \mathbf{g}_i \mathbf{p}_i / \mathbf{p}_i^T \mathbf{H} \mathbf{p}_i \quad (9.a)$$

$$\mathbf{g}_{i+1} = \mathbf{g}_i - \alpha_i \mathbf{p}_i^T \mathbf{H} \quad (9.b)$$

$$\mathbf{p}_{i+1} = \mathbf{g}_{i+1} + \beta_i \mathbf{p}_i, \quad \beta_i = \mathbf{g}_{i+1} \mathbf{H} \mathbf{p}_i / \mathbf{p}_i^T \mathbf{H} \mathbf{p}_i \quad (9.c)$$

end

**Figure 1** Conjugate gradient algorithm for minimizing a quadratic function.

One complete loop of the standard CG algorithm is equivalent to one Newton iteration and converges in  $n$  steps for *quadratic* functions (of course, the Hessian must be positive definite). The loop must be repeated for nonquadratic functions. Orthogonality of subsequent search directions is maintained with respect to two norms: one is the inner product  $\mathbf{g}_i \mathbf{p}_i$ , the second is the ‘weighted’ inner product  $\mathbf{p}_i^T \mathbf{H} \mathbf{p}_i$ . If the computation of the Hessian is to be avoided, different estimation techniques are available for  $\alpha_i$  and  $\beta_i$ , usually at the cost of slower convergence. The step size  $\alpha_i$  is the solution to  $\min\{E(\mathbf{x}_i + \alpha \mathbf{p}_i)\}$  and can be approximated using a line search. For  $\beta_i$  well known estimation techniques are available [3][7].

### 4.3 Scaled Conjugate Gradient Algorithm

The scaled conjugate gradient (SCG) algorithm circumvents the multiplication  $\mathbf{H}\mathbf{p}$  and computes the transpose of this vector directly, without computing the Hessian. An elegant way to compute the *exact* product of a Hessian and an arbitrary vector has been independently rediscovered for neural networks in [7][8][10]. This method is used for the numerical inversion proposed in this paper, using the *exact* values for  $\alpha_i$  and  $\beta_i$ . We begin with the first-order expansion of the error gradient,

$$\mathbf{g}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{g}(\mathbf{x}) + \Delta\mathbf{x}^T \mathbf{H}. \quad (10)$$

Setting  $\Delta\mathbf{x} = r\mathbf{p}$  in Equation 10, with scalar  $r$  and vector  $\mathbf{p}$ , gives

$$r\mathbf{p}^T \mathbf{H} = \mathbf{g}(\mathbf{x} + r\mathbf{p}) - \mathbf{g}(\mathbf{x}). \quad (11)$$

The desired product is obtained by dividing Equation 11 by  $r$  and taking the limit,

$$\mathbf{p}^T \mathbf{H} = \lim_{r \rightarrow 0} \frac{\mathbf{g}(\mathbf{x} + r\mathbf{p}) - \mathbf{g}(\mathbf{x})}{r} = \frac{\partial}{\partial r} [\mathbf{g}(\mathbf{x} + r\mathbf{p})]_{r=0} = R_p \{ \mathbf{g}(\mathbf{x}) \}, \quad (12)$$

which coincides with the definition of a derivative. The definition of the differential operator  $R_p \{ \cdot \}$  is based on the work in [8]. The usual rules for differential operators apply. In this paper, we obtain the desired product by applying this operator to those equations which compute gradients. In the following we simplify the notation to  $R \{ \cdot \}$ .

## 5 Using Feedforward Neural Networks

The above second order optimization methods are applied to the numerical inversion of nonlinear functions represented by the MLP in Equation 1. The error function is the usual squared error,

$$E(\mathbf{x}) = \frac{1}{2} \cdot \mathbf{e}^T(\mathbf{x}) \cdot \mathbf{e}(\mathbf{x}) = \frac{1}{2} \cdot [\mathbf{f}(\mathbf{x}) - \mathbf{y}_d]^T \cdot [\mathbf{f}(\mathbf{x}) - \mathbf{y}_d]. \quad (13)$$

Minimizing  $E$  solves the nonlinear equation represented by the network  $f$ , i.e. it finds the network input  $\mathbf{x}$  for a given desired output  $\mathbf{y}_d$ .

### 5.1 Derivation of the Algorithm

The application of Newton's method and the standard CG algorithm is straightforward if the Hessian of  $E$  is known. We apply the SCG algorithm for computing  $\mathbf{H}\mathbf{p}$  [5]. With  $\mathbf{g} = E'$ ,  $\mathbf{f}' = \mathcal{J}/\partial\mathbf{x}$ ,

$$\mathbf{p}^T \mathbf{H} = R\{\mathbf{g}(\mathbf{x})\} = R\{\mathbf{e}^T(\mathbf{x}) \cdot \mathbf{f}'(\mathbf{x})\} = R\{\mathbf{e}^T(\mathbf{x})\} \cdot \mathbf{f}'(\mathbf{x}) + \mathbf{e}^T(\mathbf{x}) \cdot R\{\mathbf{f}'(\mathbf{x})\}. \quad (14)$$

The two  $R\{\cdot\}$  terms in Equation 14 are (with  $\mathbf{a} = \mathbf{V}\mathbf{x}$ ,  $\mathbf{a}_p = \mathbf{V}(\mathbf{x} + r\mathbf{p})$ ,  $\mathbf{c} = \mathbf{V}\mathbf{p}$ ):

$$R\{\mathbf{e}^T(\mathbf{x})\} = (\mathbf{W} \cdot \Sigma'(\mathbf{a}) \cdot \mathbf{V} \cdot \mathbf{p})^T, \quad (15)$$

$$R\{\mathbf{f}'(\mathbf{x})\} = R\{\mathbf{W} \cdot \Sigma'(\mathbf{a}) \cdot \mathbf{V}\} = \mathbf{W} \cdot R\{\Sigma'(\mathbf{a})\} \cdot \mathbf{V} \quad (16)$$

$$= \left( \left[ \frac{\partial \mathbf{a}_p^T \mathbf{V}^T}{\partial r} \right]_{r=0} \otimes \mathbf{I}_{\text{row}(\Sigma')} \right) \cdot \left( \mathbf{I}_{\text{col}(r)} \otimes \left[ \frac{\partial \Sigma'(\mathbf{a}_p)}{\partial \mathbf{a}_p} \right]_{r=0} \right) \quad (17)$$

$$= (\mathbf{c}^T \otimes \mathbf{I}_q) \cdot \left( \mathbf{1} \otimes \frac{\partial \Sigma'(\mathbf{a})}{\partial \mathbf{a}} \right) = \begin{bmatrix} \mathbf{c}^T & & \\ & \ddots & \\ & & \mathbf{c}^T \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial \Sigma'(\mathbf{a})}{\partial a_1} \\ \vdots \\ \frac{\partial \Sigma'(\mathbf{a})}{\partial a_q} \end{bmatrix} \quad (18)$$

$$= \begin{bmatrix} c_1 \cdot \sigma''(a_1) & & \\ & \ddots & \\ & & c_q \cdot \sigma''(a_q) \end{bmatrix} = \Sigma_c''(\mathbf{a}). \quad (19)$$

The operator  $\otimes$  performs the Kronecker tensor product [1]. Using the above derivations gives

$$\mathbf{p}^T \mathbf{H} = [\mathbf{p}^T \cdot \mathbf{V}^T \cdot \Sigma'(\mathbf{a}) \cdot \mathbf{W}^T \cdot \mathbf{W} \cdot \Sigma'(\mathbf{a}) \cdot \mathbf{V}] + [\mathbf{e}^T(\mathbf{x}) \cdot \mathbf{W} \cdot \Sigma_c''(\mathbf{a}) \cdot \mathbf{V}], \quad (20)$$

whose software implementation can easily be optimized for computational speed. Thus we have derived a form of Equation 12 based on the specific structure of an MLP defined by Equation 1. Consequently, if such an MLP is trained to represent a given application problem we need only go into the trained MLP and read out the values of  $\mathbf{V}$ ,  $\mathbf{W}$ ,  $\mathbf{b}$  for plugging into Equation 20 and then into Equation 9.

### 5.2 Example

We assume an application context in which the process to be solved is being represented by a MLP (Equation 1) which was trained using data from the process. We use the second order methods presented above (Newton, CG, SCG) to invert the nonlinear function which is represented by the neural network. In our example we use a two-variable system, and after training our (two input, two hidden PE, two output) MLP has the following parameter values and desired output  $\mathbf{y}_d$ :

$$\mathbf{V} = \begin{bmatrix} 1 & 0.4 \\ 0.6 & 1.5 \end{bmatrix}, \mathbf{W} = \begin{bmatrix} 1 & 0.5 \\ 0.4 & 1.0 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}, \mathbf{y}_d = \begin{bmatrix} 1.0440 \\ 0.8867 \end{bmatrix}. \quad (21)$$

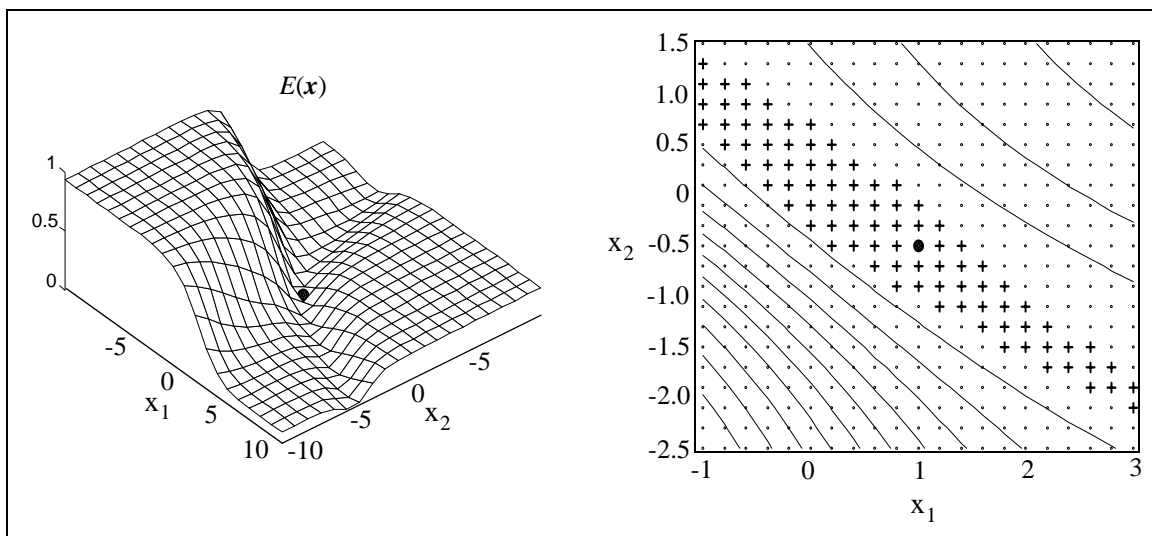
Since we constructed this example for demonstration purposes, we already know the desired fixed point  $\mathbf{x}^* = [1, -0.5]$  and choose  $\mathbf{x}_1 = [0, -0.2]$  as a convenient initial condition within the positive def-

inite neighborhood of  $\mathbf{x}^*$ , as required by the theory. We choose  $\|e\|_\infty \leq 0.00001$  as the convergence criterion for the algorithms. From a mathematical perspective all three second-order methods compute precisely the same trajectory, but the efficiency of their implementations varies. Newton's method (Equation 8) is the most complex, because it not only computes the Hessian of the error function to be minimized, but also inverts it at each iteration. The standard CG algorithm (Figure 1) avoids the matrix inversion, but still needs the Hessian itself. The most efficient second-order technique of the three is the SCG algorithm which uses Equation 20 in the algorithm of Figure 1, to avoid both, computation and inversion of the Hessian. As shown in the previous section, Equation 12 takes the form of Equation 20 for the MLP in Equation 1. The values of Equation 21 were used in the experiment tabulated in Table 1. As is well known (and demonstrated here), the second-order techniques offer dramatic improvement in the number of iterations to converge over the first order (gradient descent) technique. The 'learning rate' was 1. It will be noticed there is similarly dramatic improvement in the amount of (total) computation required to achieve the solution. This was the objective of the algorithm we are presenting here. Although this example represents a small scale problem, we obtained similar results with larger networks.

For our example problem, we know the underlying equation, so we can determine the shape of the  $E(\mathbf{x})$  surface (Figure 2). The minimum (our desired solution point) is shown as a black dot. The second figure 'zooms in' on the neighborhood of the minimum and shows the definiteness of the surface. Only a small fraction is positive definite (marked with '+'). All second-order techniques require that the initial estimate (starting point) must be in this region. This is a limitation of the methodology.

**Table 1: Algorithm Performances in the Example.**

	SCG Algorithm	CG Algorithm	Newton's Method	Gradient Descent
Iterations	7	7	7	974
Flops	1924	2141	2582	53620
Fixed Point	[0.9999,-0.5000]	[0.9999,-0.5000]	[0.9999,-0.5000]	[0.9987,-0.4991]



**Figure 2** Error surface and definiteness of the 2-2-2 multilayer perceptron used in the example.

### 5.3 Applicable Class of Problems

Important assumptions must be made when using second order optimization methods. First, the non-linear function is (at least locally) invertible, i.e. input and output space have the same dimension. Second, a good estimate of the desired fixed point is required for choosing an initial condition in the positive definite neighborhood of that point. Third, the optimization process must not leave the positive definite portion of the error surface, otherwise the algorithm immediately can become unstable. This is a stringent constraint, since in the usual case only a small portion of the error surface is positive definite, as is shown in Figure 2 for our example. If we were to randomly select an initial condition within this 'mountainous' area, there is no reason to expect that it will be in that positive definite region.

Techniques have been proposed to overcome some of these difficulties [3][7], but any additional feature attached to the algorithms will cost computation time, whereas our goal is *efficient* problem solving for (possibly) time-critical applications. We rather impose constraints on the class of problems intended to be solved, and therefore assume *a priori* knowledge about the topology of the error function. This is reasonable near minima, where most functions behave like quadratics. Consider a scenario where the minimum is initially known (as in the example) and is used as the starting point. If the application context changes such that its corresponding error surface moves away from that point a sufficiently small step, then the point is not the minimum any longer, but is still in the positive definite neighborhood of the new minimum. It can therefore be used as the new initial condition. If this process is repeated, the subsequent optimization processes 'chase' the moving minimum. An example of this class of problems is the inverse kinematics problem in robotics, where a sequence of joint angles is needed to guide the robot's end effector along a desired Cartesian trajectory. The minimum is initially known - one simply measures the joint angles. The robot arm can then be moved in small steps, while an optimization process continuously finds the best joint angles, never leaving the positive definite neighborhood of the (moving) minimum.

## 6 Conclusion

An efficient second-order optimization algorithm for the numerical inversion of nonlinear functions represented by feedforward neural networks was designed and its performance and limitations evaluated. As is well known, second-order information can dramatically reduce the number of iterations for convergence. It was demonstrated, that the algorithm developed herein, based on a neural network, also dramatically reduces the amount of computations required. Stringent conditions on the application context must be satisfied for a successful application of these algorithms. The method was demonstrated with an example, and a class of applicable problems was briefly characterized.

## 7 References

- [1] J.W. Brewer, "Kronecker Products and Matrix Calculus in System Theory," *IEEE Transactions on Circuits and Systems*, Vol. 25, No. 9, pp. 772-781, September 1978.
- [2] J. Dieudonné, *Foundations of Modern Analysis*, Academic Press, New York, 1960.
- [3] M.R. Hestenes, *Conjugate Direction Methods in Optimization*, Springer-Verlag, New York, 1980.
- [4] K. Hornik, M. Stinchcombe and H. White, "Multilayer Feedforward Networks are Universal Approximators," *Neural Networks*, Vol. 2, pp. 359-366, 1989.
- [5] K. Mathia, *Solutions of Linear Equations and a Class of Nonlinear Equations Using Recurrent Neural Networks*, Ph.D. Dissertation, Portland State University, Portland, Oregon, 1996.
- [6] K. Mathia and R. Saeks, "Solving Nonlinear Equations Using Recurrent Neural Networks," *Proc. World Congress on Neural Networks 1995*, Washington D.C., Vol. 1, pp. 76-79, July 1995.
- [7] M.F. Møller, "A Scaled Conjugate Gradient Algorithm for Fast Supervised Learning," *Neural Networks*, Vol. 6, No. 4, pp. 525-533, 1993.
- [8] B.A. Pearlmutter, "Fast Exact Multiplication by the Hessian," *Neural Computation*, Vol. 6, No. 1, pp. 147-160, January 1994.
- [9] P.J. Werbos, "Beyond regression: New Tools for Prediction and Analysis in the Behavioral Sciences," Ph.D. Thesis, Harvard University, Cambridge, MA, 1974.
- [10] P.J. Werbos, "Backpropagation: Past and Future," *Proc. IEEE Int. Conf. Neural Networks*, Vol. 1, pp. 343-353, San Diego, CA, 1988.