

D.C. Park, M. El-Sharkawi and R.J. Marks II, "An adaptively trained neural network", IEEE Transactions on Neural Networks, vol.2, pp.334-345, (1991).

An Adaptively Trained Neural Network

Dong C. Park, *Member, IEEE*, Mohamed A. El-Sharkawi, *Senior Member, IEEE*,
and Robert J. Marks II, *Senior Member, IEEE*

Abstract—A training procedure is proposed that adapts the weights of a trained layered perceptron artificial neural network to training data originating from a slowly varying nonstationary process. The resulting adaptively trained neural network (ATNN), based on nonlinear programming techniques, is shown to adapt to new training data that is in conflict with earlier training data without affecting the neural networks' response to data elsewhere. The adaptive training procedure also allows for new data to be weighted in terms of its significance. The adaptive algorithm is applied to the problem of electric load forecasting and is shown to significantly outperform the conventionally trained layered perceptron.

I. INTRODUCTION

IN the training of a layered perceptron, an assumption of stationarity of the training data is typically made. In a number of cases of interest, however, the training data constitute a slowly varying nonstationary process. In this paper, we present a training procedure applicable to such cases. In order for the layered perceptron's weights to adapt to a slowly varying nonstationarity, such a procedure should 1) still respond appropriately to previous training data if those data are not in conflict with the new training data and 2) adapt to the new training data even when they are in conflict with portions of the old data.

We propose a procedure for such adaptation which is applicable when the training data's stationarity varies sufficiently slowly. Our procedure for adaptive updating ensures a proper response to previous training data by seeking to minimize a weight sensitivity cost function while minimizing the mean square error normally ascribed to the layered perceptron. The process is illustrated through application to an exemplar interpolation problem and through its use on an electric load forecasting problem with data generated by the power industry.

II. FORMULATION OF THE PROBLEM

Assume a layered perceptron artificial neural network trained with N sets of data,

$$\{(x(1), d(1)), (x(2), d(2)), \dots, (x(N), d(N))\}$$

where $x(i)$ and $d(i)$ represent the input and desired output for the i th data set and $1 \leq i \leq N$. We assume that $x(i)$

is an I -dimensional vector and $d(i)$ is scalar. The layered perceptron is assumed to have one hidden layer with h hidden neurons. The case for more than one hidden layer is given in Appendix I. The matrix W represents the weight matrix between the input and hidden neurons and v denotes the weight vector which links the hidden and output neurons. The dimensions of W and v are $I \times h$ and $h \times 1$, respectively.

For a given input data vector, $x(i)$, the output of the layered perceptron of $y(i)$ is given by

$$y(i) = f[v^T u] \quad (1)$$

$$u = f[W^T x(i)] \quad (2)$$

where $u = [u_1, u_2, \dots, u_h]^T$; u_j , $1 \leq j \leq h$, represents the activation of the j th hidden neuron; the superscript T denotes the transpose of a matrix or vector; $f[\cdot]$ is the sigmoid function; $f[x] = 1/(1 + \exp(-x))$, $x \in \mathcal{R}$; and $f[b]$ is the $h \times 1$ vector function

$$f[b] = [f_1[b], f_2[b], \dots, f_h[b]]^T.$$

The sigmoid function for each of the hidden neurons is assumed to be identical, $f_1[\cdot] = f_2[\cdot] = \dots = f_h[\cdot] = f[\cdot]$. We assume that $W(N)$ and $v(N)$ are the weights that minimize the error function [1]:

$$E(N) = \frac{1}{2} \sum_{i=1}^N (d(i) - y(i))^2. \quad (3)$$

III. PROBLEM STATEMENT

The statement of our problem can now be described as follows:

Given $W(N)$, $v(N)$, the N sets of data, and $(x(N+1), d(N+1))$, determine $W(N+1)$ and $v(N+1)$ such that

$$\begin{aligned} E(N+1) &= \frac{1}{2} \sum_{i=1}^{N+1} (d(i) - y(i))^2 \\ &= E(N) + \frac{1}{2} (d(N+1) - y(N+1))^2 \end{aligned} \quad (4)$$

is minimized in such a manner that $y(N+1) \approx d(N+1)$.

IV. LINEARIZATION PROCESS

Linearization of the problem allows for analytic and implementation tractability: For $(x(N+1), d(N+1))$, (1) and (2) yield

$$d(N+1) = f[v^T(N+1)u] = f[(v(N) + \Delta v)^T u] \quad (5)$$

Manuscript received May 7, 1990; revised December 12, 1990. This work was supported by the Puget Sound Power and Light Company and by the Washington Technology Center at the University of Washington.

D. C. Park is with the Department of Electrical and Computer Engineering, Florida International University, Miami, FL 33199.

M. A. El-Sharkawi and R. J. Marks II are with the Department of Electrical Engineering, University of Washington, Seattle, WA 98195.

IEEE Log Number 9143023.

where

$$\begin{aligned} \mathbf{u} &= \mathbf{f}[\mathbf{W}^T(N+1) \mathbf{x}(N+1)] \\ &= \mathbf{f}[(\mathbf{W}(N) + \Delta\mathbf{W})^T \mathbf{x}(N+1)]. \end{aligned} \quad (6)$$

We expand (5) and (6) in a truncated Taylor series about the state of interest, $\{\mathbf{x}(N+1), \mathbf{W}(N+1)\}$, that is closest to the current operating state, $\{\mathbf{x}(N), \mathbf{W}(N)\}$. Such a linearization is used in other adaptive filter techniques, such as the extended Kalman filtering [2], [3] and quasi linearization [4]–[6].

Let

$$\mathbf{b} = \mathbf{W}^T(N) \mathbf{x}(N+1)$$

and

$$\Delta\mathbf{b} = \Delta\mathbf{W}^T \mathbf{x}(N+1),$$

Application of a first-order Taylor series expansion to (6) gives

$$\begin{aligned} \mathbf{u} &= \mathbf{f}[\mathbf{W}^T(N) \mathbf{x}(N+1) + \Delta\mathbf{W}^T \mathbf{x}(N+1)] \\ &= \mathbf{f}[\mathbf{b} + \Delta\mathbf{b}] \\ &\approx \mathbf{f}[\mathbf{b}] + (\nabla_{\mathbf{b}} \mathbf{f}[\mathbf{b}]) \Delta\mathbf{b} \end{aligned}$$

where $\nabla_{\mathbf{b}} \mathbf{f}[\mathbf{b}]$ is the gradient of $\mathbf{f}[\mathbf{b}]$ with respect to \mathbf{b} with elements $\partial f_i[\mathbf{b}]/\partial b_j$. This approximation is valid when $H \ll 1$, where H is the Hessian of $\mathbf{f}[\mathbf{b}]$. We require that the perturbation, $\Delta\mathbf{b}$ (or $\Delta\mathbf{W}$), be small enough so that

$$\Delta\mathbf{b} \ll \mathbf{f}[\mathbf{b}]/H$$

where the inequality applies to each component of the vector.

Since $f_i[\cdot] = f[\cdot]$, the i th component of $\mathbf{f}[\mathbf{b}]$ is only a function of b_i (that is, $f_i[\mathbf{b}] = f_i[b_i]$, $1 \leq i \leq h$) and

$$\begin{aligned} \partial f_i[\mathbf{b}]/\partial b_j &= \partial f[b_i]/\partial b_j \\ &= f[b_i] (1 - f[b_i]) \delta_{i-j} \\ &= u_i^* (1 - u_i^*) \delta_{i-j} \end{aligned}$$

then,

$$\begin{aligned} \nabla_{\mathbf{b}} \mathbf{u}^* &= \text{diag} [u_1^*(1 - u_1^*), u_2^*(1 - u_2^*), \dots, \\ &\quad u_h^*(1 - u_h^*)] \end{aligned}$$

where δ_k , the Kronecker delta, is 1 for $k = 0$ and is zero otherwise. Also, $u_i^* = f[b_i]$ is the activation of the i th hidden neuron for the new input data with the old weight such as

$$\mathbf{u}^* = [u_1^*, u_2^*, \dots, u_h^*]^T = \mathbf{f}[\mathbf{W}^T(N) \mathbf{x}(N+1)].$$

Therefore, the activations of hidden neurons are given by

$$\mathbf{u} \approx \mathbf{u}^* + (\nabla_{\mathbf{b}} \mathbf{u}^*) \Delta\mathbf{W}^T \mathbf{x}(N+1). \quad (7)$$

Inverting (5) yields

$$f^{-1}[d(N+1)] = \mathbf{v}^T(N) \mathbf{u} + \Delta\mathbf{v}^T \mathbf{u} \quad (8)$$

where $f^{-1}[x] = \ln(x/(1-x))$.

From (7) and (8) we find that

$$\begin{aligned} f^{-1}[d(N+1)] - \mathbf{v}^T(N) \mathbf{u}^* \\ \approx \mathbf{v}^T(N) (\nabla_{\mathbf{b}} \mathbf{u}^*) \Delta\mathbf{W}^T \mathbf{x}(N+1) \\ + \Delta\mathbf{v}^T \mathbf{u}^* + \Delta\mathbf{v}^T (\nabla_{\mathbf{b}} \mathbf{u}^*) \Delta\mathbf{W}^T \mathbf{x}(N+1). \end{aligned} \quad (9)$$

Remember that the perturbation in $\Delta\mathbf{W}$ is assumed to be small in order to use (7). If $\Delta\mathbf{v} \ll \mathbf{v}$, then the third term of the right-hand side of (9) can be ignored and

$$\begin{aligned} f^{-1}[d(N+1)] - \mathbf{v}^T(N) \mathbf{u}^* \\ \approx \mathbf{v}^T(N) (\nabla_{\mathbf{b}} \mathbf{u}^*) \Delta\mathbf{W}^T \mathbf{x}(N+1) + \Delta\mathbf{v}^T \mathbf{u}^*. \end{aligned} \quad (10)$$

We rearrange $\Delta\mathbf{W}$ into vector form as follows:

$$\Delta\mathbf{W}_{\text{vec}} = [\Delta\mathbf{w}_1^T \cdots \Delta\mathbf{w}_p^T]^T = [\Delta\mathbf{W}_{\text{vec},1} \cdots \Delta\mathbf{W}_{\text{vec},p}]$$

where $\Delta\mathbf{w}_i$ is the i th row in $\Delta\mathbf{W}$, and $p = h \times I$. Note that p is the number of interconnections between the neurons in the input and the hidden layers. Then (10) can be rewritten as

$$\begin{aligned} c_1 &= [\Delta\mathbf{W}_{\text{vec}}^T; \Delta\mathbf{v}^T] \begin{bmatrix} \mathbf{u}^* \\ \cdot \\ \mathbf{u}^* \end{bmatrix} \\ &= \mathbf{z}^T \mathbf{a} \end{aligned} \quad (11)$$

where c_1 , \mathbf{a} , and \mathbf{z} are vectors defined as

$$c_1 \equiv f^{-1}[d(N+1)] - \mathbf{v}^T(N) \mathbf{u}^* \quad (12)$$

$$\mathbf{a} \equiv [\mathbf{u}^{*T}; \mathbf{u}^{*T}]^T \quad (13)$$

$$\mathbf{z} \equiv [\Delta\mathbf{W}_{\text{vec}}^T; \Delta\mathbf{v}^T]^T \quad (14)$$

and \mathbf{u}^{\dagger} is a solution of

$$\Delta\mathbf{W}_{\text{vec}}^T \mathbf{u}^{\dagger} = \mathbf{v}^T(N) \mathbf{Q} \Delta\mathbf{W}^T \mathbf{x}(N+1). \quad (15)$$

Since there are $h \times (I+1)$ unknowns with one equation, (11), there exist many solutions. Uniqueness is imposed by the additional requirement that our solution be constrained by the additional conditions in the problem statement. Specifically, we need to change the weights for the $(N+1)$ st datum with minimum effect on the previous N data. We are therefore motivated to find the sensitivity of $y(i)$ over a weight change. Equation (4) can be rewritten as

$$E(N) = \frac{1}{2} \sum_{i=1}^N (d(i) - y(i))^2 = \sum_{i=1}^N E_i.$$

The sensitivity for the input weights follows as

$$\frac{\partial E_i}{\partial w_{jk}} = -[d(i) - y(i)] \left(\frac{\partial y(i)}{\partial w_{jk}} \right) \quad (16)$$

and those for the output weights are

$$\frac{\partial E_i}{\partial v_k} = -[d(i) - y(i)] \left(\frac{\partial y(i)}{\partial v_k} \right) \quad (17)$$

where w_{jk} , the weight of interconnection between input neuron j and hidden neuron k , is the jk th element of \mathbf{W} .

Let

$$y(i) = f[\text{sum}_y] \quad \text{and} \quad u_k = f[\text{sum}_u]$$

where

$$\text{sum}_y = \sum_k u_k v_k \quad \text{and} \quad \text{sum}_u = \sum_j w_{jk} x_j.$$

Thus,

$$\begin{aligned} \frac{\partial y(i)}{\partial v_k} &= \left(\frac{\partial y(i)}{\partial \text{sum}_y} \right) \left(\frac{\partial \text{sum}_y}{\partial v_k} \right) = \left. \frac{\partial f[x]}{\partial x} \right|_{x=y(i)} u_k \\ &= y(i) (1 - y(i)) u_k \\ &\equiv S V_{i,k} \end{aligned} \quad (18)$$

and

$$\begin{aligned} \frac{\partial y(i)}{\partial w_{jk}} &= \left(\frac{\partial y(i)}{\partial u_k} \right) \left(\frac{\partial u_k}{\partial w_{jk}} \right) \\ &= \left(\frac{\partial y(i)}{\partial \text{sum}_y} \right) \left(\frac{\partial \text{sum}_y}{\partial u_k} \right) \left(\frac{\partial u_k}{\partial \text{sum}_u} \right) \left(\frac{\partial \text{sum}_u}{\partial w_{jk}} \right) \\ &= \left. \frac{\partial f[x]}{\partial x} \right|_{x=y(i)} v_k \left. \frac{\partial f[x]}{\partial x} \right|_{x=u_k} x_j \\ &= y(i) (1 - y(i)) v_k u_k (1 - u_k) x_j \\ &\equiv S W_{i,jk} \end{aligned} \quad (19)$$

where $S V_{i,k}$ is the sensitivity of $y(i)$ caused by small changes in v_k and $S W_{i,jk}$ is the sensitivity of $y(i)$ to the $w_{j,k}$'s. Consequently,

$$\Delta E_i = \sum_{j,k} \left(\frac{\partial E_i}{\partial w_{jk}} \right) \Delta W_{jk} + \sum_k \left(\frac{\partial E_i}{\partial v_k} \right) \Delta v_k. \quad (20)$$

By combining (18), (19), and (20) with (16) and (17), we obtain

$$\begin{aligned} \Delta E_i &= (\epsilon_i) [\Delta W_{\text{vec},1} \cdots \Delta W_{\text{vec},p} \Delta v_1 \cdots \Delta v_h] \\ &\quad \cdot [S W_{i,1} \cdots S W_{i,p} \cdots S V_{i,1} \cdots S V_{i,h}]^T \end{aligned} \quad (21)$$

where $\epsilon_i = -(d(i) - y(i))$, $1 \leq i \leq N$, and $p = I \times h$. Equivalently, using (20), we can write the matrix equation

$$\Delta \mathbf{E} = \mathbf{A} \mathbf{S} \mathbf{z} \quad (22)$$

where the i th element of $\Delta \mathbf{E}$ is ΔE_i :

$$\mathbf{A} = \text{diag} [\epsilon_1, \cdots, \epsilon_N] \quad (23)$$

and

$$\mathbf{S} = \begin{bmatrix} S W_{1,1} \cdots S W_{1,p} & S V_{1,1} \cdots S V_{1,h} \\ \vdots & \vdots \\ S W_{N,1} \cdots S W_{N,p} & S V_{N,1} \cdots S V_{N,h} \end{bmatrix} \quad (24)$$

where \mathbf{z} is a $q \times 1$ vector defined in (11) and $q = p + h = (I + 1) \times h$.

The weights $\mathbf{W}(N)$ and $\mathbf{v}(N)$ are optimal for minimizing $E(N)$. With the addition of the $(N + 1)$ st datum, $d(N + 1) = y(N + 1)$, the objective function of (4) can be changed to

$$\begin{aligned} J &= \frac{1}{2} \sum_{i=1}^N (E_{i,W(N)} - E_{i,W(N+1)})^2 \\ &= \frac{1}{2} \sum_{i=1}^N \Delta E_i^2 \end{aligned}$$

where $E_{i,W(N)}$ and $E_{i,W(N+1)}$ are the errors of the i th datum with $\{\mathbf{W}(N), \mathbf{v}(N)\}$ and $\{\mathbf{W}(N + 1), \mathbf{v}(N + 1)\}$, respectively. Equivalently, from (22),

$$\begin{aligned} J &= \frac{1}{2} (\Delta \mathbf{E})^T (\Delta \mathbf{E}) \\ &= \frac{1}{2} (\mathbf{A} \mathbf{S} \mathbf{z})^T (\mathbf{A} \mathbf{S} \mathbf{z}) \\ &= \frac{1}{2} \mathbf{z}^T \mathbf{K} \mathbf{z} \end{aligned} \quad (25)$$

where $\mathbf{K} = \mathbf{S}^T (\mathbf{A}^T \mathbf{A}) \mathbf{S}$. Note that $\mathbf{K} = \mathbf{K}^T$.

Since there exists only one equation with q unknowns in (11), the solutions of (11) are on a $(q - 1)$ -dimensional hyperplane. Since small perturbations for \mathbf{W} and \mathbf{v} are assumed in (7), among the solutions in (11), only those with small variation in weight space (i.e., small $\|\mathbf{z}\|$) are allowable.

The problem is now a standard nonlinear programming problem. Specifically,

$$\begin{aligned} \text{minimize} \quad & J(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{K} \mathbf{z} \\ \text{subject to} \quad & \mathbf{z}^T \mathbf{a} = c_1. \end{aligned}$$

A equicost line of the cost function $J(\mathbf{z})$ represents a q -dimensional hyperellipse centered at the origin of the \mathbf{z} plane. The constraint $\mathbf{z}^T \mathbf{a} = c_1$ is a $(q - 1)$ -dimensional hyperplane on the \mathbf{z} plane. The shape of the cost function is dictated by the eigenvalues and eigenvectors of \mathbf{K} [7], [8].

When there is no boundary constraint or too loose of a boundary constraint for the value of \mathbf{z} , the choice of \mathbf{z} (say \mathbf{z}^*) may minimize the cost index $J(\mathbf{z})$, at the cost of severely violating the linearity assumption. If we introduce a boundary constraint, \mathcal{B} , for \mathbf{z} , we can avoid this problem at a cost of being at a low though not minimum value of the cost function. However, if we use too restrictive (or tight) a boundary, as shown in Fig. 1, a solution that satisfies the two different constraints may not exist.

The above discussion suggests that we need to limit the perturbations to make them small enough to meet the linearization assumption of the truncated Taylor series. At the same time, we need to allow \mathbf{z} to be large enough to have at least one solution on the $\mathbf{z}^T \mathbf{a} = c_1$ plane.

V. BOUNDARY CONSTRAINT

Define the boundary as a set of points in the \mathbf{z} plane:

$$\mathcal{B} = \{\mathbf{z}: -\bar{\ell} \leq \mathbf{z} \leq \bar{\ell}\}$$

where $\bar{\ell} = [l_1, l_2, \cdots, l_q]^T$ and $l_i \geq 0$, $i = 1, \cdots, q$.

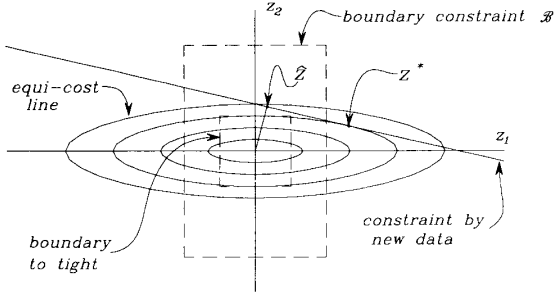


Fig. 1. Contours of constant $J(z)$ when z is of dimension 2. The constraint imposed by the new data is shown as a line.

The $(q - 1)$ -dimensional constraint hyperplane formed by $z^T a = c_1$ contains a unique point, \hat{z} , that is minimal distance from the origin. The boundary constraint should include at least \hat{z} . In other words, we need to choose \bar{l} that includes at least \hat{z} . For the plane $z^T a = c_1$, the result is

$$\hat{z} = \frac{c_1 a}{a^T a}. \quad (26)$$

By using \hat{z} given in (26), we choose \bar{l} such as $\bar{l} = k\hat{z}$, $k \geq 1$. This boundary ensures intersection by the linear constraint plane. A minimal distance point and an example of a properly selected boundary constraint for the two-dimensional z are shown in Fig. 1. The larger the k , the smaller $J(z)$ can be. Choice of an overly large value for k , however, may violate the Taylor series approximation used in (7).

VI. REDUCED GRADIENT METHOD

Using the boundary constraint, our problem now becomes

$$\begin{aligned} &\text{minimize} && J(z) = \frac{1}{2} z^T K z \\ &\text{subject to} && z^T a = c_1 \\ &&& \text{and } z \in \mathcal{B}, \text{ or } -\bar{l} \leq z \leq \bar{l}. \end{aligned} \quad (27)$$

There exists a variety of nonlinear programming methods to solve this type of problem (e.g., the gradient projection method and the convex simplex method [9]–[11]). We found the *reduced gradient method* (RGM) using a linear constraint most conducive to the problem.

The RGM starts with the partitioning of variables into two groups: basic and nonbasic. The variables in the basic group are dependent and the variables in the nonbasic group are independent. The cardinality of the basic group is the number of linear constraints. In our case, we will have one basic variable since there is one given linear constraint. The other $(q - 1)$ variables are assigned as independent variables.

Assume that we have a feasible solution, z , which satisfies the constraints in (27). We partition the q -dimensional variable z arbitrarily into two groups: $z = (z_\alpha, z_\beta) = [z_\alpha, z_\beta^T]^T$, where z_α is a one-dimensional basic

(dependent) variable and z_β is a $(q - 1)$ -dimensional nonbasic (independent) variable. The problem in (27) can then be written as

$$\text{minimize} \quad J(z_\alpha, z_\beta) = \frac{1}{2} (z_\alpha, z_\beta)^T K (z_\alpha, z_\beta) \quad (28)$$

$$\text{subject to} \quad z_\alpha a_\alpha + a_\beta^T z_\beta = c_1 \quad (29)$$

$$-\bar{l} \leq (z_\alpha, z_\beta) \leq \bar{l} \quad (30)$$

where $a = [a_\alpha, a_\beta^T]^T$.

The RGM minimizes the objective function iteratively only in terms of independent variables. The movement of a dependent variable is dictated by the relationship in (29). That is, if $z_\beta(k + 1) = z_\beta(k) + \Delta z_\beta$, and $z_\alpha(k + 1) = z_\alpha(k) + \Delta z_\alpha$, where k is an iteration index, then

$$\Delta z_\alpha a_\alpha + a_\beta^T \Delta z_\beta = 0$$

or

$$\Delta z_\alpha = -\frac{1}{a_\alpha} a_\beta^T \Delta z_\beta. \quad (31)$$

The movement controlled by (31) guarantees that the new point, $(z_\alpha(k + 1), z_\beta(k + 1))$, is always on the $(q - 1)$ -dimensional hyperplane given in (29).

The steepest descent method is now used to find the direction of the independent variable in the RGM. The gradient of each variable determines the direction of movement at each step of the steepest descent method. The gradient of $J(z)$ follows as

$$\nabla_z J(z) = (\nabla_{z_\alpha} J(z), \nabla_{z_\beta} J(z)) = Kz = K(z_\alpha, z_\beta).$$

Since the relationship between the movements of z_α and z_β should satisfy (31), the amount of movement, $\nabla_{z_\alpha} J(z)$, for z_α can be decomposed into the amount of movement for z_β . The gradient with respect to z_β (the *reduced gradient* of z), instead of z , is found to be

$$r^T = \nabla_{z_\beta} J(z_\alpha, z_\beta) - \frac{1}{a_\alpha} \nabla_{z_\alpha} J(z_\alpha, z_\beta) a_\beta \quad (32)$$

where $r = [r_1, r_2, \dots, r_{q-1}]^T$.

Since the boundary is a box, each variable has a corresponding boundary defined by two sides of the box. In order to allow the possibility that an independent variable is about to violate the boundary constraint, we introduce the *working set* concept. A constraint is defined as *active* if the corresponding variable reaches its constraint boundary which it is about to violate. At each step of the iteration, a subset of boundary constraints that are *active* at the current step is chosen. We call this the working set, denoted by $\mathcal{W}(z_\beta)$.

The vector z_β moves in the direction of its gradient unless it violates its boundary constraint. That is, at each step, the direction of z_β 's movement is

$$\Delta z_{\beta,i} = \begin{cases} -r_i, & i \notin \mathcal{W}(z_\beta) \\ 0, & i \in \mathcal{W}(z_\beta). \end{cases}$$

If an independent variable reaches an edge of its boundary constraint, it becomes part of the working set. This pro-

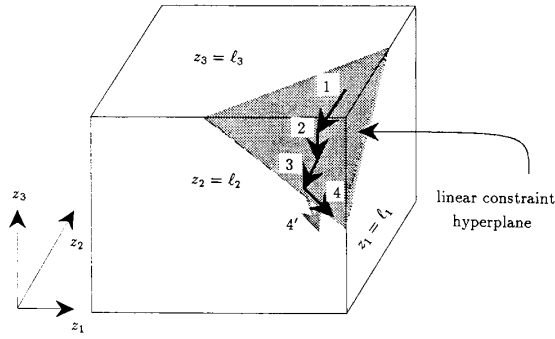


Fig. 2. Boundary constraints and searching.

cess is illustrated in Fig. 2. The numbers in Fig. 2 denote the resultant minimum cost points from the RGM. During iteration, intermediate solutions move as $0 \rightarrow 1 \rightarrow 2 \rightarrow 3$ and the boundary of z_2 is reached. This occurs, even though the gradient information implies use of path $4'$ to minimize the given cost function. The intermediate solution should not choose the path $4'$ since it violates the boundary constraint. Instead, by choosing path 4, the solution guarantees $z_2 \leq l_2$. In this case, the constraint $z_2 \leq l_2$ is called active and the corresponding variable z_2 is included in the working set. A more detailed explanation can be found in [9].

Assume that a point is found to be such that

$$r_i = 0 \quad \forall i \notin \mathcal{W}(z_\beta)$$

but there exists $j \in \mathcal{W}(z_\beta)$ such that either $r_j < 0$ and j have been put into the working set because the corresponding variable violated its lower boundary constraint, or $r_j > 0$ and j have been put into the working set because the corresponding variable violated its upper boundary constraint. Then j is deleted from the working set. Once the amount of movement, Δz_β , of z_β is found, Δz_α is given by (31).

The RGM generally converges to a relative minimum in most of the cases. There also exists a modified RGM which yields global convergence [9]. The convex simplex method is another variation of RGM which guarantees global convergence [12]. We found little performance variation in our simulations.

We desire the solution by the RGM to be the global minimum within the boundary constraint. This is ensured if the boundary constraint \mathcal{B} is a convex set and the cost function, $J(z)$, is a convex function. Both conditions are present in our problem. We investigate the globalness of the relative minimum obtained with RGM by utilizing the convexity properties of the cost function, $J(z)$, and boundary constraint, \mathcal{B} . The following definitions establish the convexity of a function and a set.

Definition 1 (Convex Set): A set Ω in E^n is said to be convex if $\forall x_1, x_2 \in \Omega$ and for $0 \leq \alpha \leq 1$, the point

$$x_3 = \alpha x_1 + (1 - \alpha)x_2$$

is also in Ω . Clearly, since \mathcal{B} is a box, it is convex.

TABLE I
A SUMMARY ATNN ALGORITHM

0. Begin
1. Calculate A and S by using N data as described in (23) and (24).
2. Find K such that $K = S^T(A^T A)S$.
3. By using (13)–(14), find the linear constraint equation:

$$z^T a = c_1.$$
4. Find the boundary constraint, $\vec{l} = 2\hat{z}$, by using

$$\hat{z} = \frac{c_1 a}{a^T a}.$$
5. Perform the reduced gradient method and find z which minimizes the cost function:

$$J(z) = \frac{1}{2} z^T K z$$
 with the constraints $z^T a = c_1$ and $-\vec{l} \leq z \leq \vec{l}$.
6. Find ΔW and Δv by rearranging z and resulting ΔW_{vec} .

$$[\Delta W_{\text{vec}}^T; \Delta v^T] = z^T.$$
7. Update $W(N)$ and $v(N)$:

$$W(N+1) = W(N) + \Delta W$$

$$v(N+1) = v(N) + \Delta v.$$
8. Stop.

Definition 2 (Convex Function): A function f defined on convex set Ω is said to be convex if $\forall x_1, x_2 \in \Omega$ and $\forall \alpha \in \mathbb{R}, 0 \leq \alpha \leq 1$,

$$f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2).$$

Lemma 1: The cost function $J(z)$ is convex. (A proof of Lemma 1 is given in Appendix II).

Since $J(z)$ is a convex function defined on the convex set Ω , any achieved relative minimum of $J(z)$ is the global minimum.

A summary of the adaptively trained ATNN is given in Tables I and II.

VII. SUSCEPTIBILITY TO DRIFT

A recursive algorithm for adaptive training must possess two features: it should adapt to new data and it should not drift. Drift is an undesired feature, wherein the weights of the layered perceptron change even when the new data do not change the input/output relationship of the existing neural network.

Assume that the trained network with $\{W(N), v(N)\}$ has the following functional relationship between the input, x , and the output, y :

$$y = f(W(N), v(N); x) \quad (33)$$

and

$$D = \{(x(i), d(i)); 1 \leq i \leq N\}.$$

If we use additional data, $(x(N+1), d(N+1))$, that are not in D but satisfy (33), both ΔW and Δv should be equal to $\mathbf{0}$. Otherwise, the proposed adaptive procedure is said to drift. We will now show that our proposed algorithm does not drift.

From (11) and (13),

$$y = d(N+1) = f(W(N), v(N); x(N+1))$$

TABLE II
A SUMMARY OF THE REDUCED GRADIENT METHOD

-
0. Set initial feasible solution $\mathbf{z} = \hat{\mathbf{z}}$.
 1. Set $z_\alpha = z_1$ and $\mathbf{z}_\beta = [z_2, z_3, \dots, z_q]^T$. Initialize $\mathcal{W}(\mathbf{z}_\beta) = \{\phi\}$, where $\{\phi\}$ denotes the empty set.
 2. Calculate

$$\mathbf{r}^T = \nabla_{\mathbf{z}_\beta} J(z_\alpha, \mathbf{z}_\beta) - \frac{1}{a_\alpha} \nabla_{z_\alpha} J(z_\alpha, \mathbf{z}_\beta) a_\beta.$$
 3. Find the direction of the movement of \mathbf{z}_β :

$$\Delta \mathbf{z}_{\beta,i} = \begin{cases} -r_i, & i \notin \mathcal{W}(\mathbf{z}_\beta) \\ 0, & i \in \mathcal{W}(\mathbf{z}_\beta). \end{cases}$$
 4. Reconstruct the working set $\mathcal{W}(\mathbf{z}_\beta)$ if necessary.
 5. We stop if the point is an acceptable solution. This occurs when $\|\Delta \mathbf{z}\| \leq \epsilon$, where ϵ is a convergence measure, or $\forall i \in \mathcal{W}(\mathbf{z}_\beta)$. Otherwise, find

$$\Delta z_\alpha = -\frac{1}{a_\alpha} \mathbf{a}_\beta^T \Delta \mathbf{z}_\beta.$$
 6. Find γ_1, γ_2 , and γ_3 such that

$$\begin{aligned} \max \{ \gamma_1: -l_\alpha \leq z_\alpha + \gamma_1 \Delta z_\alpha \leq l_\alpha; \gamma_1 \geq 0 \} \\ \max \{ \gamma_2: -\bar{l}_\beta \leq \mathbf{z}_\beta + \gamma_2 \Delta \mathbf{z}_\beta \leq \bar{l}_\beta; \gamma_2 \geq 0 \} \\ \min \{ \gamma_3: J(\mathbf{z} + \gamma_3 \mathbf{z}); 0 \leq \gamma_3 \leq \gamma_1, 0 \leq \gamma_3 \leq \gamma_2 \}. \end{aligned}$$
 7. Calculate $\mathbf{z} = \mathbf{z} + \gamma_3 \mathbf{z}$.
 8. If $\gamma_3 < \gamma_1$, go to 2. Otherwise, declare the dependent variable to be independent and declare one of the independent variables, which is positively inside of the nonlinearity constraint, to be dependent. Update a_α and \mathbf{a}_β . Go to 2.
-

and we have

$$c_1 = 0.$$

Consequently, $\mathbf{z}^T \mathbf{a} = 0$ becomes the line constraint for this problem. From (26), the closest point which satisfies this line constraint is

$$\hat{\mathbf{z}} = \frac{c_1 \mathbf{a}}{\mathbf{a}^T \mathbf{a}} = \mathbf{0}.$$

Since $\hat{\mathbf{z}} = \mathbf{0}$ gives

$$\mathbf{J} = \frac{1}{2} \mathbf{z}^T \mathbf{K} \mathbf{z} = 0$$

the initial point to the RGM, $\hat{\mathbf{z}} = \mathbf{0}$, is the solution. This implies that

$$\Delta \mathbf{W} = \mathbf{0} \quad \text{and} \quad \Delta \mathbf{v} = \mathbf{0}.$$

Therefore, we conclude that there is no drifting in the ATNN. Note, alternatively, that the drift problem can also be addressed by triggering the ATNN only when the output error exceeds some specified threshold.

VIII. CHOOSING WEIGHTS ON THE E_i 's

The weighting coefficient, γ_i , on each of the error terms of (3) can be used to weight the error function as follows:

$$\hat{E}(N) = \frac{1}{2} \sum_{i=1}^N \gamma_i (d(i) - y(i))^2 = \sum_{i=1}^N \hat{E}_i$$

where $\hat{E}(N)$ is a weighted version of $E(N)$ and

$$\hat{E}_i = \frac{1}{2} \gamma_i (d(i) - y(i))^2.$$

By introducing this weighting coefficient, $\gamma_i = \mathcal{F}(i)$, we can control the comparative importance of the previous data.

The weighting coefficients can be incorporated into the ATNN as follows. Since

$$\hat{E}_i = \frac{1}{2} \gamma_i (d(i) - y(i))^2 = \gamma_i E_i$$

we can rewrite (22) with the WCGF:

$$\Delta \hat{E} = \mathbf{\Gamma} \mathbf{A} \mathbf{S} \mathbf{z} = \hat{\mathbf{A}} \mathbf{S} \mathbf{z}$$

where $\hat{\mathbf{A}} = \mathbf{\Gamma} \mathbf{A}$ and

$$\mathbf{\Gamma} = \text{diag} [\gamma_1, \gamma_2, \dots, \gamma_N].$$

Therefore, we can consider the weighting coefficients in the ATNN by simply changing \mathbf{A} to $\hat{\mathbf{A}}$.

IX. EXPERIMENTS AND RESULTS

Since Lapedes and Farber [14] suggested their use for the system identification, the layered perceptron has been successfully applied to many other system identification problems [14]–[23]. In the following example ATNN applications, the layered perceptron is used for system identification.

First, an exemplar problem that has single input and single output is considered. It nicely illustrates the proposed algorithm graphically. Then, the ATNN is applied to the electric load forecasting problem in order to illustrate its application to highly nonstationary training data.

A. Exemplar Problem

We are given a set of one-dimensional input and one-dimensional output data such that

$$\{(x(1), d(1)), (x(2), d(2)), \dots, (x(100), d(100))\}.$$

where $x(n) = 0.01(n - 1)$. A 3-layered perceptron with ten hidden neurons was trained using error back-propagation to obtain $\mathbf{W}(100)$ and $\mathbf{v}(100)$, the weights for the given 100 data points.

In order to illustrate the ATNN, assume that a new data point $(x(101), d(101))$ is to be used to update the network weights and that the new data are

$$x(101) = x(51) = 0.5$$

and

$$d(101) = 1.1 \times d(51) = 0.45 + 0.045 = 0.495.$$

Thus, the new data are inconsistent with some of previous training data.

Results: We applied our algorithm to the above prob-

TABLE III
THE RESULTANT CHANGE IN WEIGHTS FROM USING THE ATNN

	Old	Δ	New		Old	Δ	New
$W_{0,2}$	-1.6941	0.0274	-1.6667	$W_{1,2}$	-2.3778	-0.0006	-2.3785
$W_{0,3}$	-1.2054	0.0514	-1.1539	$W_{1,3}$	-0.5910	-0.0549	-0.6459
$W_{0,4}$	-1.1859	0.2209	-0.9649	$W_{1,4}$	-7.5534	-0.2781	-7.2752
$W_{0,5}$	15.3013	-0.4345	14.8667	$W_{1,5}$	-22.4426	0.1881	-22.2545
$W_{0,6}$	-1.0096	0.0190	-0.9906	$W_{1,6}$	-0.9404	-0.0181	-0.9586
$W_{0,7}$	-1.1323	0.0375	-1.0947	$W_{1,7}$	-0.4236	-0.0484	-0.4720
$W_{0,8}$	-1.0451	0.0260	-1.0191	$W_{1,8}$	-0.4856	-0.0337	-0.5193
$W_{0,9}$	-2.0645	0.1131	-1.9514	$W_{1,9}$	-3.9084	0.0350	-3.8733
$W_{0,10}$	2.1602	0.0125	2.1727	$W_{1,10}$	-14.4980	-0.2369	-14.7349
$W_{1,11}$	-0.9274	0.0330	-0.8944	$W_{1,11}$	-0.5305	-0.0427	-0.5732
v_0	-1.4489	-0.3237	-1.7726	v_1	-0.6528	0.6081	-0.0446
v_2	-0.1127	-0.1584	-0.2712	v_3	-0.4850	-0.1392	-0.6242
v_4	-4.1520	-0.0143	-4.1376	v_5	2.0000	0.3647	2.3647
v_6	-0.0708	-0.1916	-0.2625	v_7	-0.3882	-0.1280	-0.5163
v_8	-0.2301	-0.1413	-0.3715	v_9	-1.1492	-0.0908	-1.2400
v_{10}	1.9850	0.4237	2.4087	v_{11}	-0.2795	-0.1582	-0.4377

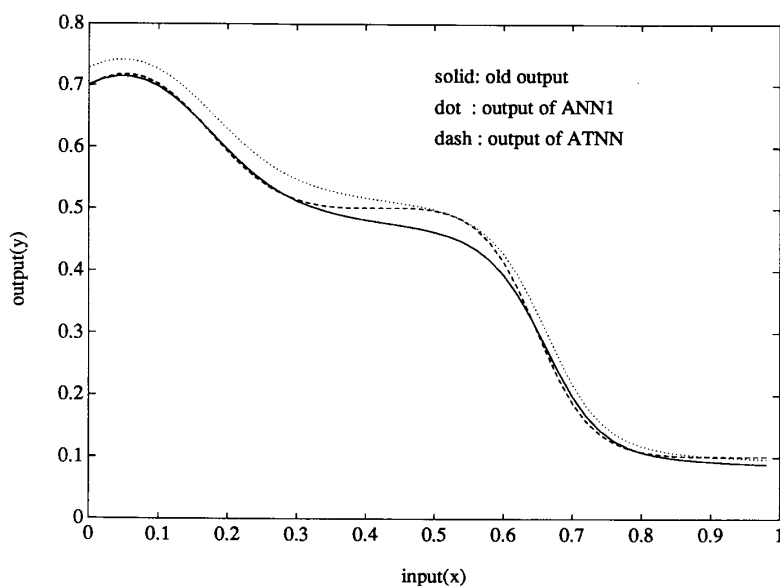


Fig. 3. A performance comparison of the ATNN (dashed line) to (a) previously trained layered perceptron (solid line) and (b) a layered perceptron retained with new data (dotted line). The performance of a layered perceptron retained with both old and new data is graphically indistinguishable from (a).

lem and found a set of optimal weights, $W(101)$ and $v(101)$, as shown in Table III. In order to show the effects of the new weights on the previous data, we test the layered perceptron with the new weights. Fig. 3 shows the changes in outputs for the previous data. Note that the layered perceptron responds to the desired output of the new data despite inconsistency between the old and the new data. (The representation obtained from the neural net trained with the old data is graphically indistinguishable from the training data.)

A comparison is also made between the ATNN and two other training approaches. In the first alternative method,

a layered perceptron is retrained with only the new data point using the old network as a point of initialization. This guarantees that the network, after retraining, produces a desirable output for the new data. As illustrated in Fig. 3, the retrained network, however, does not maintain proper performance for the old data. Yet another layered perceptron was trained with both the old and new data. The layered perceptron trained with the old data was used for the initialization. The retrained network with this technique, however, has a performance that is graphically indistinguishable from that of the performance of the layered perceptron trained only on the old data.

B. Application to Load Forecasting

Electric load forecasting is used by electric utility companies for prediction of electric power consumption based on pertinent data such as past consumption and weather information.

Electric system load forecasting with lead times from hours to several days can assist a power system operator to efficiently schedule spinning reserve allocation. In addition, the ability to forecast electric system load can also provide valuable information for possible energy interchange with other utilities. Accurate information for future loads is also useful in detecting many vulnerable situations in advance if applied to system security assessment problem. Because of the importance of load forecasting, forecasts should be accurate and sufficiently timely to assist the system operator.

1) *Approaches*: Previous approaches to load forecasting in general fall into three categories. First, the load pattern is treated as a time series signal and future loads are predicted by using various time series analysis techniques [24]–[26]. The second approach assumes that the load pattern is highly correlated with such weather variables as temperature and wind speed and finds a functional relationship between weather variables and the system load. The future load can then be forecast by inserting the predicted weather information into the predetermined functional relationship [27]–[30]. The last approach combines the time series approach and the regression approach by using artificial neural networks. More details on this application and corresponding results can be found in [22], where forecasting was performed by training a layered perceptron with conventional error back-propagation. A comparative study on different problems, including the load forecasting problem, contrasting the performance of a layered perceptron and the classification and regression trees (CART) has been recently reported by Atlas *et al.* [23].

2) *Experiments*: Hourly temperature and load data for the Seattle/Tacoma area from November 1, 1989, to February 9, 1990, were collected and provided for us by the Puget Sound Power and Light Company. Our focus is on hourly load forecasting with lead time of 48 hours for a normal weekday (i.e., no holidays or weekends). All data except the last nine days (February 1–9) are used for training the layered perceptron, and the data of the last nine days are used for testing the trained layered perceptron. As a topology for the layered perceptron, a good choice of input variables for the output, L_k , was found to be

$$(k, T_{k-168}, T_{k-50}, T_{k-49}, T_{k-48}, L_{k-168}, L_{k-50}, L_{k-49}, L_{k-48}, T_k)$$

where k is the hour of the day and T_i and L_j represents the temperature and load at hour i and j respectively. The layered perceptron used in this experiment has 11 inputs (including a bias term), five hidden neurons, and one output neuron.

TABLE IV
A PERFORMANCE COMPARISON OF THE ATNN TO REGULAR LAYERED PERCEPTRON BY THE AVERAGE OF ABSOLUTE ERROR (%) IN THE ELECTRIC LOAD FORECASTING PROBLEM

	2/1	2/2	2/7	2/8	2/9
ANN	6.14	6.53	5.28	2.92	2.36
ATNN	2.18	2.26	4.65	2.36	2.79

First, a standard error back-propagation algorithm is used for training the layered perceptron. The ATNN algorithm is then applied to the trained layered perceptron to improve the performance. The ATNN adapts the trained layered perceptron by sequentially using the latest 24 data points while the regular layered perceptron simply augments the training data set with 24 new entries.

3) *Results*: The actual temperatures are used in the testing stage to find the effects of the ATNN algorithm on the forecast loads since forecast temperatures deviate from actual temperature by as much as 14% (February 1, 1990). Use of forecast temperatures therefore distracts the ATNN.

Daily average errors are calculated as

$$\text{percent error} = \frac{1}{24} \sum_{k=1}^{24} \frac{|L_k - L'_k|}{L_k} \times 100\%$$

where L_k and L'_k denote actual and forecast loads at hour k .

As shown in Table IV, the comparative results imply that the ATNN gives improved results except for day 5. The improvement is especially significant on day 2, when the weather changed abruptly. Fig. 4 depicts a comparative example of five days of forecasts.

These results were generated using a value of $k = 2$. As with the step size and momentum terms in error back-propagation, the choice of a good value of k remains somewhat of an art.

X. CONCLUSION

In this paper, an adaptive training algorithm for a layered perceptron has been proposed. The algorithm uses a previously trained layered perceptron to adapt the weights for the new data. Conventional learning algorithms; such as back-propagation, can give an averaged result for inconsistent data, which can occur often in a nonstationary environment.

This adaptively trained neural network (ATNN) utilizes nonlinear programming methods to adapt to temporally nonstationary training data.

We have demonstrated the use of the ATNN by applying it to the nonstationary problem of electric load forecasting, where the seasonal or annual load change can be severe. The ATNN has the additional flexibility of allowing importance weighting on training data.

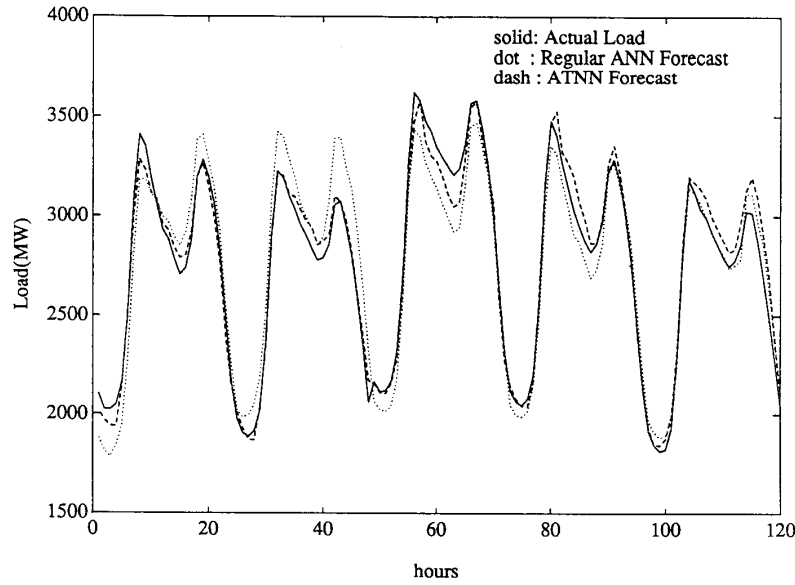


Fig. 4. A performance comparison of the ATNN (dashed line) to conventional layered perceptron (dotted line) when applied to electric load forecasting problem. Solid line represents actual load.

APPENDIX I SENSITIVITY MATRIX FOR THE GENERAL CASE

Consider a layered perceptron which has n hidden layers and each hidden layer has N_i hidden neurons. Assume that $w_{i,j}^{(m)}$ is the weight of connection between neuron i in the m th hidden layer from the output layer and the neuron j in the $(m-1)$ st hidden layer from the output layer. Then, $\mathbf{u}^{(i)} = [u_1^{(i)}, u_2^{(i)}, \dots, u_{N_i}^{(i)}]^T$, where N_i represents the number of neurons in the i th hidden layer from the output layer and $u_j^{(i)}$ denotes the activation of the neuron j in that hidden layer. Then the sensitivity of the ℓ th output, y_ℓ , by the change of $w_{i,j}^{(m)}$ is given by the following.

1) For $m = 1$:

$$\frac{\partial y_\ell}{\partial w_{i,\ell}^{(m)}} = \frac{\partial f \left[\sum_{j=1}^{N_i} w_{j,\ell}^{(1)} \right]}{\partial w_{i,\ell}^{(1)}} = y_\ell (1 - y_\ell) u_i^{(1)}.$$

2) $m = 2$:

$$\frac{\partial y_\ell}{\partial w_{i,j}^{(m)}} = \left(\frac{\partial y_\ell}{\partial w_j^{(1)}} \right) \left(\frac{\partial u_j^{(1)}}{\partial u_i^{(2)}} \right).$$

3) For $m \geq 3$:

$$\frac{\partial y_\ell}{\partial w_{i,j}^{(m)}} = \nabla_{\mathbf{u}^{(1)}} y_\ell \left\{ \prod_{k=1}^{m-3} \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right) \right\} \left(\frac{\partial \mathbf{u}^{(m-2)}}{\partial \mathbf{u}_j^{(m-1)}} \right) \left(\frac{\partial u_j^{(m-1)}}{\partial u_i^{(m)}} \right) \quad \text{and} \quad (34)$$

where

$$u_p^{(k)} = f \left[\sum_{q=1}^{N_{(k+1)}} w_{q,p}^{(k+1)} u_q^{(k+1)} \right]$$

$$\begin{aligned} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} &= \frac{\partial [u_1^{(k)}, u_2^{(k)}, \dots, u_{N_k}^{(k)}]^T}{\partial \mathbf{u}^{(k+1)}} \\ &= \begin{bmatrix} \nabla_{\mathbf{u}^{(k+1)}} u_1^{(k)} \\ \vdots \\ \nabla_{\mathbf{u}^{(k+1)}} u_{N_k}^{(k)} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial u_1^{(k)}}{\partial u_1^{(k+1)}} & \dots & \frac{\partial u_1^{(k)}}{\partial u_{N_{(k+1)}}^{(k+1)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial u_{N_k}^{(k)}}{\partial u_1^{(k+1)}} & \dots & \frac{\partial u_{N_k}^{(k)}}{\partial u_{N_{(k+1)}}^{(k+1)}} \end{bmatrix} \\ &= \begin{bmatrix} \delta_1^{(k)} w_{1,1}^{(k+1)} & \dots & \delta_1^{(k)} w_{N_{(k+1)},1}^{(k+1)} \\ \vdots & \ddots & \vdots \\ \delta_{N_k}^{(k)} w_{1,N_k}^{(k+1)} & \dots & \delta_{N_k}^{(k)} w_{N_{(k+1)},N_k}^{(k+1)} \end{bmatrix} \end{aligned}$$

$$\delta_i = \left. \frac{\partial f[x]}{\partial x} \right|_{x=u_i^{(k)}} = u_i^{(k)} (1 - u_i^{(k)}).$$

Proof: For the cases where $m > 3$, define

$$\frac{\partial u_j^{(k)}}{\partial u_i^{(k+1)}} \stackrel{\text{def}}{=} \gamma(i(k+1), j(k)).$$

Since

$$\frac{\partial \mathbf{u}^{(k-1)}}{\partial \mathbf{u}^{(k+1)}} = \left(\frac{\partial \mathbf{u}^{(k-1)}}{\partial \mathbf{u}^{(k)}} \right) \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right)$$

the sensitivity of $u_n^{(k-1)}$ to the change of $u_i^{(k+1)}$ is given by

$$\begin{aligned} \frac{\partial u_n^{(k-1)}}{\partial u_i^{(k+1)}} &= \left(\frac{\partial u_n^{(k-1)}}{\partial \mathbf{u}^{(k)}} \right) \left(\frac{\partial \mathbf{u}^{(k)}}{\partial u_i^{(k+1)}} \right) \\ &= \nabla_{\mathbf{u}^{(k)}} u_n^{(k-1)} \begin{bmatrix} \frac{\partial u_1^{(k)}}{\partial u_i^{(k+1)}} \\ \vdots \\ \frac{\partial u_{N_m}^{(k)}}{\partial u_i^{(k+1)}} \end{bmatrix} \\ &= \nabla_{\mathbf{u}^{(k)}} u_n^{(k-1)} \begin{bmatrix} \gamma(i(k+1), 1(k)) \\ \vdots \\ \gamma(i(k+1), N_k(k)) \end{bmatrix}. \end{aligned}$$

Thus, the sensitivities of activation in the $(k-1)$ st hidden layer by a change in $u_i^{(k+1)}$ are

$$\begin{aligned} \frac{\partial \mathbf{u}^{(k-1)}}{\partial u_i^{(k+1)}} &= \frac{\partial}{\partial u_i^{(k+1)}} [u_1^{(k-1)}, \dots, u_{N_{(k-1)}}^{(k-1)}]^T \\ &= \left[\frac{\partial u_1^{(k-1)}}{\partial u_i^{(k+1)}}, \dots, \frac{\partial u_{N_{(k-1)}}^{(k-1)}}{\partial u_i^{(k+1)}} \right]^T \\ &= \begin{bmatrix} \nabla_{\mathbf{u}^{(k)}} u_1^{(k-1)} \begin{bmatrix} \gamma(i(k+1), 1(k)) \\ \vdots \\ \gamma(i(k+1), N_k(k)) \end{bmatrix} \\ \vdots \\ \nabla_{\mathbf{u}^{(k)}} u_{N_{(k-1)}}^{(k-1)} \begin{bmatrix} \gamma(i(k+1), 1(k)) \\ \vdots \\ \gamma(i(k+1), N_k(k)) \end{bmatrix} \end{bmatrix} \\ &= \begin{bmatrix} \nabla_{\mathbf{u}^{(k)}} u_1^{(k-1)} \\ \vdots \\ \nabla_{\mathbf{u}^{(k)}} u_{N_{(k-1)}}^{(k-1)} \end{bmatrix} \begin{bmatrix} \gamma(i(k+1), 1(k)) \\ \vdots \\ \gamma(i(k+1), N_k(k)) \end{bmatrix} \\ &= \frac{\partial [u_1^{(k-1)}, \dots, u_{N_{(k-1)}}^{(k-1)}]^T}{\partial \mathbf{u}^{(k)}} \\ &= \begin{bmatrix} \gamma(i(k+1), 1(k)) \\ \vdots \\ \gamma(i(k+1), N_k(k)) \end{bmatrix} \end{aligned}$$

As a result, the sensitivities of activations for the neurons in the $(k-1)$ st hidden layer by the change of activations for the neurons in the $(k+1)$ st hidden layer are

$$\begin{aligned} \frac{\partial \mathbf{u}^{(k-1)}}{\partial \mathbf{u}^{(k+1)}} &= \left[\frac{\partial u_1^{(k-1)}}{\partial u_1^{(k+1)}}, \dots, \frac{\partial u_{N_{(k-1)}}^{(k-1)}}{\partial u_{N_{(k-1)}}^{(k+1)}} \right] \\ &= \frac{\partial [u_1^{(k-1)}, \dots, u_{N_{(k-1)}}^{(k-1)}]^T}{\partial \mathbf{u}^{(k)}} \\ &= \begin{bmatrix} \gamma(1(k+1), 1(k)) & \cdots & \gamma(N_{(k+1)}(k+1), 1(k)) \\ \vdots & \ddots & \vdots \\ \gamma(1(k+1), N_k(k)) & \cdots & \gamma(N_{(k+1)}(k+1), N_k(k)) \end{bmatrix} \\ &= \frac{\partial \mathbf{u}^{(k-1)}}{\partial \mathbf{u}^{(k)}} \begin{bmatrix} \frac{\partial u_1^{(k)}}{\partial u_1^{(k+1)}} \\ \vdots \\ \frac{\partial u_{N_k}^{(k)}}{\partial u_{N_k}^{(k+1)}} \end{bmatrix} \\ &= \left(\frac{\partial \mathbf{u}^{(k-1)}}{\partial \mathbf{u}^{(k)}} \right) \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right). \end{aligned}$$

Therefore, the sensitivities of activation for the neurons from one layer by the changes of activations for the neurons to the other layer can be calculated in a cascaded fashion. That is,

$$\begin{aligned} \frac{\partial \mathbf{u}^{(0)}}{\partial \mathbf{u}^{(m)}} &= \prod_{k=0}^{m-1} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \\ &= \frac{\partial \mathbf{u}^{(0)}}{\partial \mathbf{u}^{(1)}} \left\{ \prod_{k=1}^{m-3} \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right) \right\} \\ &\quad \cdot \left(\frac{\partial \mathbf{u}^{(m-2)}}{\partial \mathbf{u}^{(m-1)}} \right) \left(\frac{\partial \mathbf{u}^{(m-1)}}{\partial \mathbf{u}^{(m)}} \right) \end{aligned} \quad (35)$$

where $\mathbf{u}^{(0)}$ are the activations of neurons in output layer.

Since the change of $w_{i,j}^{(m)}$ affects only $u_j^{(m-1)}$ by an amount $\partial u_j^{(m-1)} / \partial w_{i,j}^{(m)}$, the change of $u_j^{(m-1)}$ affects all the activations of neurons in the $(m-1)$ st hidden layer by the amount of $\partial \mathbf{u}^{(m-2)} / \partial u_j^{(m-1)}$, the change of $\mathbf{u}^{(m-2)}$ affects the neurons in the first hidden layer by $\partial \mathbf{u}^{(1)} / \partial \mathbf{u}^{(m-2)}$, and, finally, the change of $\mathbf{u}^{(1)}$ affects the output neuron y_ℓ by $\partial y_\ell / \partial \mathbf{u}^{(1)}$. Therefore,

$$\begin{aligned} \frac{\partial y_\ell}{\partial w_{i,j}^{(m)}} &= \frac{\partial y_\ell}{\partial \mathbf{u}^{(1)}} \left\{ \prod_{k=1}^{m-3} \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right) \right\} \left(\frac{\partial \mathbf{u}^{(m-2)}}{\partial u_j^{(m-1)}} \right) \left(\frac{\partial u_j^{(m-1)}}{\partial u_i^{(m)}} \right) \\ &= \nabla_{\mathbf{u}^{(1)}} y_\ell \left\{ \prod_{k=1}^{m-3} \left(\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{u}^{(k+1)}} \right) \right\} \left(\frac{\partial \mathbf{u}^{(m-2)}}{\partial u_j^{(m-1)}} \right) \left(\frac{\partial u_j^{(m-1)}}{\partial u_i^{(m)}} \right). \end{aligned}$$

This concludes the proof of (34). Q.E.D.

By using the sensitivity of each weight to the output, the error resulting from the perturbed weights can be found as in (20):

$$\Delta E_i = \sum_{j,k,m} \left(\frac{\partial E_i}{\partial w_{j,k}^{(m)}} \right) \Delta w_{j,k}^{(m)}.$$

The sensitivity matrix S for the general case can also be found by the same method applied for the single hidden layer case by using (22) to (24).

APPENDIX II

PROOF OF LEMMA 1 (CONVEXITY OF THE COST FUNCTION)

For our case,

$$J(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{K} \mathbf{z}.$$

Thus

$$\begin{aligned} J(\alpha \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2) &= (\alpha \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2)^T \mathbf{K} (\alpha \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2) \\ &= \alpha \mathbf{z}_1^T \mathbf{K} \alpha \mathbf{z}_1 + \alpha \mathbf{z}_1^T \mathbf{K} (1 - \alpha) \mathbf{z}_2 \\ &\quad + (1 - \alpha) \mathbf{z}_2^T \mathbf{K} \alpha \mathbf{z}_1 + (1 - \alpha)^2 \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2 \\ &= \alpha^2 \mathbf{z}_1^T \mathbf{K} \mathbf{z}_1 + \alpha(1 - \alpha) [\mathbf{z}_1^T \mathbf{K} \mathbf{z}_2 + \mathbf{z}_2^T \mathbf{K} \mathbf{z}_1] \\ &\quad + (1 - \alpha)^2 \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2. \end{aligned} \quad (36)$$

Let

$$\bar{\alpha} = \alpha J(\mathbf{z}_1) + (1 - \alpha) J(\mathbf{z}_2)$$

then

$$\bar{\alpha} = \alpha \mathbf{z}_1^T \mathbf{K} \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2.$$

Also,

$$\begin{aligned} \bar{\alpha} &= \alpha \bar{\alpha} + (1 - \alpha) \bar{\alpha} \\ &= \alpha^2 \mathbf{z}_1^T \mathbf{K} \mathbf{z}_1 + \alpha(1 - \alpha) \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2 \\ &\quad + (1 - \alpha) \alpha \mathbf{z}_1^T \mathbf{K} \mathbf{z}_1 + (1 - \alpha)^2 \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2. \end{aligned} \quad (37)$$

By combining (36) and (37),

$$\begin{aligned} \bar{\alpha} - J(\alpha \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2) &= \alpha(1 - \alpha) [\mathbf{z}_1^T \mathbf{K} \mathbf{z}_1 + \mathbf{z}_2^T \mathbf{K} \mathbf{z}_2 - \mathbf{z}_1^T \mathbf{K} \mathbf{z}_2 - \mathbf{z}_2^T \mathbf{K} \mathbf{z}_1] \\ &= \alpha(1 - \alpha) (\mathbf{z}_1 - \mathbf{z}_2)^T \mathbf{K} (\mathbf{z}_1 - \mathbf{z}_2) \geq 0. \end{aligned}$$

Therefore, $\alpha J(\mathbf{z}_1) + (1 - \alpha) J(\mathbf{z}_2) \geq J(\alpha \mathbf{z}_1 + (1 - \alpha) \mathbf{z}_2)$ and consequently $J(\mathbf{z})$ is a convex function by Definition 2. Q.E.D.

ACKNOWLEDGMENT

The authors thank M. L. Bruce of the Puget Sound Power and Light Company for providing the data for the electric load forecasting problem and Dr. S. Oh, M. Damborg, J.-N. Hwang, and L. Atlas for their valuable comments. They also wish to thank Dr. H. Rauch and the reviewers for their helpful suggestions in preparing this manuscript.

REFERENCES

- [1] D. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representation by error propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, D. Rumelhart et al., Eds. Cambridge, MA: MIT Press, 1986.
- [2] A. H. Jazwinski, *Stochastic Processes and Filtering Theory*. New York: Academic Press, 1970.
- [3] A. Gelb, *Applied Optimal Estimation*. Cambridge, MA: M.I.T. Press, 1974.
- [4] R. Bellman, *Dynamic Programming*. Princeton, NJ: Princeton University Press, 1957.
- [5] K. S. P. Kumar and R. Shridar, "On the identification of control systems by the quasi-linearization method," *IEEE Trans. Automat. Contr.*, vol. AC-9, pp. 151-154, 1964.
- [6] A. P. Sage and B. R. Eisenberg, "Experiments in nonlinear and non-stationary system identification via quasilinearization and differential approximation," in *Proc. Joint Automat. Contr. Conf.*, 1965, pp. 522-530.
- [7] F. R. Gantmacher, *The Theory of Matrices*, vol. 1. New York: Chelsea Publishing, 1960.
- [8] J. N. Franklin, *Matrix Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1968.
- [9] D. Luenberger, *Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley, 1984, pp. 295-321.
- [10] A. V. Fiacco and G. P. McCormick, *Nonlinear Programming: Sequential Unconstrained Minimization Technique*. New York: Wiley, 1968.
- [11] G. P. McCormick, "Optimality criteria in nonlinear programming," in *SIAM-AMS Proc.*, IX, 1976, pp. 27-38.
- [12] W. I. Zangwill, *Nonlinear Programming: Unified Approach*. Englewood Cliffs, NJ: Prentice-Hall, 1969.
- [13] D. Graupe, *Time Series Analysis, Identification and Adaptive Filtering*. Malabar, FL: Robert E. Kreiger, 1984.

- [14] D. Lapedes and R. Farber, "Nonlinear signal processing using neural networks: prediction and system modeling," Los Alamos National Laboratory, Los Alamos, NM, TR LA-UR-87-2662, 1987.
- [15] J. Moody and D. Darken, "Learning with localized receptive fields," in *Proc. 1988 Connectionist Summer School*, 1988.
- [16] M. Tenorio and W. Lee, "Self-organizing network for optimum supervised learning," *IEEE Trans. Neural Networks*, vol. 1, Mar. 1990.
- [17] M. Aggoune, M. El-Sharkawi, D. Park, M. Damborg, and R. Marks II, "Preliminary results on using artificial neural networks for security assessment," in *Proc. 1989 Power Ind. Comp. Appl. Conf.* (Seattle, WA), May 1989. (to appear in *IEEE Trans. Power Syst.*)
- [18] M. El-Sharkawi, R. Marks II, M. Aggoune, D. Park, M. Damborg, and L. Atlas, "Dynamic security assessment of power systems using back error propagation artificial neural networks," in *Proc. 2nd Annual Symp. Expert Systems Appl. to Power Systems* (Seattle, WA), July 1989.
- [19] C. Lam *et al.*, "Determination of particle size distribution using a neural network trained with backscatter measurement," in *Proc. 1990 AP-S Int. Symp. & URSI Radio Sci. Meeting* (Dallas, TX), May 1990.
- [20] A. Ishimaru *et al.*, "Optical sensing of particle size distribution by using artificial neural networks technique," *Opt. Lett.*, vol. 15, pp. 1221-1223, 1990.
- [21] L. Atlas *et al.*, "A performance comparison of trained multilayer perceptrons and trained classification trees," in *Proc. IEEE Conf. Syst., Man, Cybern.* (Boston, MA), Nov. 1989.
- [22] D. Park, M. El-Sharkawi, R. Marks II, and L. Atlas, "Electric load forecasting using an artificial neural network," *IEEE Trans. Power Syst.*, to be published.
- [23] L. Atlas *et al.*, "A performance comparison of trained multi-layer perceptrons and trained classification trees," *Proc. IEEE*, vol. 78, pp. 1614-1619, 1990.
- [24] G. E. Box and G. M. Jenkins, *Time Series Analysis—Forecasting and Control*. San Francisco: Holden-Day, 1976.
- [25] S. Vemuri, D. Hill, and R. Balasubramanian, "Load forecasting using stochastic models," in *Proc. 8th PICA Conf.* (Minneapolis, MN), 1973, pp. 31-37.
- [26] W. Christiaanse, "Short-term load forecasting using general exponential smoothing," *IEEE Trans. Power App. Syst.*, vol. PAS-90, pp. 900-910, Apr. 1971.
- [27] P. Gupta and Y. Yamada, "Adaptive short-term forecasting of hourly loads using weather information," *IEEE Trans. Power App. Syst.*, vol. PAS-91, pp. 2085-2094, 1972.
- [28] C. Asbury, "Weather load model for electric demand energy forecasting," *IEEE Trans. Power App. Syst.*, vol. PAS-94, pp. 1111-1116, 1975.
- [29] S. Rahman and R. Bhatnagar, "An expert system based algorithm for short load forecast," *IEEE Trans. Power Syst.*, vol. 3, pp. 392-399, May 1988.
- [30] K. Jabbour, J. Riveros, D. Landbergen, and W. Meyer, "ALFA: Automated load forecasting assistant," *IEEE Trans. Power Syst.*, vol. 3, pp. 908-914, Aug. 1988.



Dong Chul Park (S'80-M'90) received the B.S. degree in 1980 from Sogang University, the M.S. degree in 1982 from the Korea Advanced Institute of Science and Technology, Seoul, Korea, and the Ph.D. degree in electrical engineering from the University of Washington, Seattle, in June 1990.

He joined the faculty of the Department of Electrical and Computer Engineering at Florida International University, the State University of Florida at Miami, in August 1990. Dr. Park was a committee member for the First International

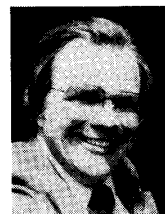
Forum on Applications of Neural Networks to Power Systems (Seattle, 1991). Since 1988 he has published 11 archival journal and 15 proceedings papers in the areas of neural networks, optical computing, signal processing, and electric load forecasting.



Mohamed A. El-Sharkawi (S'76-M'80-SM'83) was born in Cairo, Egypt, in 1948. He received the B.Sc. degree in electrical engineering in 1971 from the Cairo High Institute of Technology, Egypt. His M.A.Sc and Ph.D. degrees in electrical engineering were received from the University of British Columbia in 1977 and 1980 respectively.

Since 1980, he has been with the Department of Electrical Engineering at the University of Washington, where he is currently a Professor. He is a pioneer in the area of NN applications to power systems. He is the Chairman of the IEEE Power Engineering Society task force on applications of neural networks to power systems. He has organized and chaired several special sessions and panel discussions on the subject at various IEEE conferences. He was the organizer and chairman of the first conference to be held in Seattle, in 1990, dedicated to the applications of neural networks to power systems.

Dr. El-Sharkawi has published numerous papers in the areas of neural network applications to power systems, power system dynamics, power electronics, electric drives, and high-performance tracking and control.



Robert J. Marks II (S'76-M'77-SM'83) holds the title of Professor in the Department of Electrical Engineering at the University of Washington, Seattle. He was awarded the Outstanding Branch Councillor award in 1982 by the IEEE and, in 1984, was presented with an IEEE Centennial Medal. He was also the Chair of the IEEE Neural Networks Committee and was the cofounder and first Chair of the IEEE Circuits & Systems Society Technical Committee on Neural Systems & Applications. Dr. Marks was also elected the first President of the IEEE Council on Neural Networks. He is a Fellow of the Optical Society of America. He also was a cofounder and first President of the Puget Sound Section of the Optical Society of America and was elected that organization's first honorary member. He is cofounder and current President of the Multidimensional Systems Corporation in Bellevue, WA.

Dr. Marks is the topical editor for Optical Signal Processing and Image Science for the *Journal of the Optical Society of America—A* and is a member of the Editorial Board for *The International Journal of Neurocomputing*. He is the author of *Introduction to Shannon Sampling and Interpolation Theory* (Springer-Verlag, 1991). He has published and holds patents in the areas of signal analysis, detection theory, signal recovery, optical computing, signal processing, and artificial neural processing. Dr. Marks is a cofounder of the Christian Faculty Fellowship at the University of Washington and is a member of Eta Kappa Nu and Sigma Xi.