

# Theory of the Backpropagation Neural Network

Robert Hecht-Nielsen

HNC, Inc.  
5501 Oberlin Drive  
San Diego, CA 92121  
619-546-8877

and

Department of Electrical and Computer Engineering  
University of California at San Diego  
La Jolla, CA 92139

## Abstract

Backpropagation is currently the most widely applied neural network architecture. The information processing operation that it carries out is the approximation of a mapping or function  $f: A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , from a bounded subset  $A$  of  $n$ -dimensional Euclidean space to a bounded subset  $f[A]$  of  $m$ -dimensional Euclidean space, by means of training on examples  $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k), \dots$  of the mapping's action, where  $y_k = f(x_k)$ . It is assumed that such examples are generated by selecting  $x_k$  vectors randomly from  $A$  in accordance with a fixed probability density function  $\rho(x)$ . This paper presents a survey of the basic theory of the backpropagation neural network architecture covering the areas of: architectural design, performance measurement, function approximation capability, and learning. The survey includes previously known material, as well as some new results: a formulation of the backpropagation neural network architecture to make it a valid neural network (past formulations violated the locality of processing restriction) and a proof that the backpropagation mean squared error function exists and is differentiable. Also included is a theorem showing that any  $L_2$  function from  $[0, 1]^n$  to  $\mathbb{R}^m$  can be implemented to any desired degree of accuracy with a three-layer backpropagation neural network. Finally, an Appendix presents a speculative neurophysiological model illustrating how the backpropagation neural network architecture might plausibly be implemented in the mammalian brain for cortico-cortical learning between nearby regions of cerebral cortex.

## 1 Introduction

Without question, backpropagation is currently the most widely applied neural network architecture. This popularity primarily revolves around the ability of backpropagation networks to learn complicated multidimensional mappings. One way to look at this ability is that, in the words of Werbos [49,50,52], backpropagation goes "Beyond Regression".

Backpropagation has a colorful history. Apparently, it was originally introduced by Bryson and Ho in 1969 [5] and independently rediscovered by Werbos in 1974 [52], by Parker in the mid 1980's [41,39,40] and by Rumelhart, Williams and other members of the PDP group in 1985 [46,44,1]. Although the PDP group became aware of Parker's work shortly after their discovery (they cited Parker's 1985 report in their first papers on backpropagation [59,46]), Werbos' work was not widely appreciated until mid-1987, when it was found by Parker. The work of Bryson and Ho was discovered in 1988 by le Cun [34]. Even earlier incarnations may yet emerge.

Notwithstanding its checkered history, there is no question that credit for developing backpropagation into a usable technique, as well as promulgation of the architecture to a large audience, rests entirely with Rumelhart and the other members of the PDP group [45]. Before their work, backpropagation was unappreciated and obscure. Today, it is a mainstay of neurocomputing.

One of the crucial decisions in the design of the backpropagation architecture is the selection of a sigmoidal activation function (see Section 2 below). Historically, sigmoidal activation functions have been used by a number of investigators. Grossberg was the first advocate of the use of sigmoid functions in neural networks [18]; although his reasons for using them are not closely related to their role in backpropagation (see [8] for a discussion of the relationship between these two bodies of work). Sejnowski, Hinton, and Ackley [21,45] and Hopfield [23] provided still other reasons for using sigmoidal activation functions, but again, these are not directly related to backpropagation. The choice of sigmoid activation for backpropagation (at least for the PDP group reincarnation of the architecture) was made quite consciously by Williams, based upon his 1983 study of activation functions [60]. As we shall see, it was a propitious choice.

## 2 Backpropagation Neural Network Architecture

This Section reviews the architecture of the backpropagation neural network. The transfer function equations for each processing element are provided for both the forward and backward passes. First, we recall the definition of a neural network:

**Definition 1** *A neural network is a parallel, distributed information processing structure consisting of processing elements (which can possess a local memory and can carry out localized information processing operations) interconnected together with unidirectional signal channels called connections. Each processing element has a single output connection which branches ("fans out") into as many collateral connections as desired (each carrying the same signal - the processing element output signal). The processing element output signal can be of any mathematical type desired. All of the processing that goes on within each processing element must be completely local: i.e., it must depend only upon the current values of the input signals arriving at the processing element via impinging connections and upon values stored in the processing element's local memory.*

The importance of restating the neural network definition relates to the fact that (as pointed out by Carpenter and Grossberg [8]) traditional forms of the backpropagation architecture are, in fact, not neural networks. They violate the locality of processing restriction. The new backpropagation neural network architecture presented below eliminates this objection, while retaining the traditional mathematical form of the architecture.

The backpropagation neural network architecture is a hierarchical design consisting of fully interconnected *layers* or *rows* of processing *units* (with each unit itself comprised of several individual processing elements, as will be explained below). Backpropagation belongs to the class of *mapping neural network* architectures and therefore the information processing function that it carries out is the approximation of a bounded mapping or function  $f: A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , from a compact subset  $A$  of  $n$ -dimensional Euclidean space to a bounded subset  $f[A]$  of  $m$ -dimensional Euclidean space, by means of training on examples  $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k), \dots$  of the mapping, where  $y_k = f(x_k)$ . It will always be assumed that such examples of a mapping  $f$  are generated by selecting  $x_k$  vectors randomly from  $A$  in accordance with a fixed probability density function  $\rho(x)$ . The

operational use to which the network is to be put after training is also assumed to involve random selections of input vectors  $\mathbf{x}$  in accordance with  $p(\mathbf{x})$ . The backpropagation architecture described in this paper is the basic, classical version. Many variants of this basic form exist (see Section 5).

The macro-scale detail of the backpropagation neural network architecture is shown in Figure 1. In general, the architecture consists of  $K$  rows of processing units, numbered from the bottom up beginning with 1. For simplicity, the terms *row* and *layer* will be used interchangeably in this paper, even though each row will actually turn out to consist of two heterogeneous layers (where the term *layer* is used to denote a collection of processing elements having the same form of transfer function). The first layer consists of  $n$  fanout processing elements that simply accept the individual components  $x_i$  of the input vector  $\mathbf{x}$  and distribute them, without modification, to all of the units of the second row. Each unit on each row receives the output signal of each of the units of the row below. This continues through all of the rows of the network until the final row. The final ( $K^{\text{th}}$ ) row of the network consists of  $m$  units and produces the network's estimate  $\mathbf{y}'$  of the correct output vector  $\mathbf{y}$ . For the purposes of this paper it will always be assumed that  $K \geq 3$ . Rows 2 thru  $K - 1$  are called *hidden rows* (because they are not directly connected to the outside world).

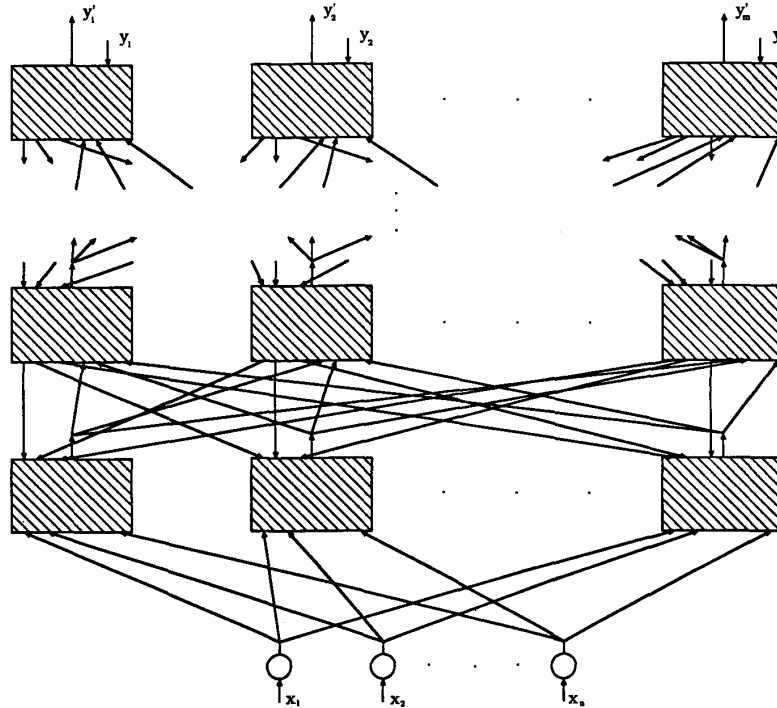


Figure 1: Macroscopic Architecture of the Backpropagation Neural Network. The boxes are called *units*. The detailed architecture of the units is elaborated in Figure 2. Each row can have any number of units desired; except the output row, which must have exactly  $m$  units (because the output vector  $\mathbf{y}'$  must be  $m$ -dimensional).

Besides the feedforward connections mentioned above, each unit of each hidden row receives an 'error feedback' connection from each of the units above it. However, as will be seen below, these are not merely fanned out copies of a broadcast output (as the forward connections are), but are each separate connections, each carrying a different signal. The details of the individual "units" (shown as rectangles in Figure 1) are revealed in Figure 2 (which depicts two units on adjacent rows and shows all of their connections). Note that each unit is composed of a single *sun* processing element and several *planet* processing elements. Each planet produces an output signal that is distributed to both its sun and to the sun of the previous layer that supplied input to it. Each planet receives input from one of the suns of the previous layer as well as its own sun. As stated above, the hidden row suns receive input from one of the planets of each of the suns on the next higher row. The output row suns receive the 'correct answer'  $y_i$  for their component of the output vector on each training trial. As discussed in detail below, the network functions in two stages: a forward pass and a backward pass. A *scheduling processing element* (not shown) sends signals to each of the processing elements of the network telling it when to apply its processing element transfer function and whether to apply the forward pass part of it or the backward pass part of it. After the transfer function is applied, the output signal is latched to the value determined during the update. This value is therefore constant until the next update. The exact equations of the processing elements of the network are given in Table 1.

The scheduling of the network's operation consists of two "sweeps" through the network. The first sweep (the *forward pass*) starts by inserting the vector  $\mathbf{x}_k$  into the network's first row, the *input* (or *fanout*) layer. The processing elements of the first layer transmit all of the components of  $\mathbf{x}_k$  to all of the units of the second row of the network. The outputs of the units of row two are then transmitted to all of the units of row three, and so on, until finally the  $m$  output units (the units of the top,  $K^{\text{th}}$ , row) emit the components of the vector  $\mathbf{y}'_k$  (the network's estimate of the desired output  $\mathbf{y}_k$ ). After the estimate  $\mathbf{y}'_k$  is emitted, each of the output units is supplied with its component of the correct output vector  $\mathbf{y}_k$ , starting the second, backward, sweep through the network (the *backward pass*). The output suns compute their  $\delta_{K_i}$ 's and transmit these to their planets. The planets then update their  $\Delta_{Kij}$  values (or update their weights if enough batching has been carried out) and then transmit the values  $w_{Kij}^{\text{old}} \delta_{K_i}$  to the suns of the previous row. The superscript "old" indicates that the (non-updated) weight value used on the forward pass is used in this calculation. This process continues until the planets of row 2 (the first hidden layer) have been updated. The cycle can then be repeated. In short, each cycle consists of the inputs to the network "bubbling up" from the bottom to the top and then the errors "percolating down" from the top to the bottom.

#### Forward Pass

**Planet  $j$  of processing element  $i$   
of row  $l$  ( $l = 2, 3, \dots, K$ ):**

Input Used:  $z_{(l-1)j}$  (where  $z_{(l-1)0} \equiv 1.0$ )  
Weight Value Used:  $w_{lij}$   
Local Memory Value Used: None  
Output:  $w_{lij} z_{(l-1)j}$   
Weight and Local Memory Value Update:  
None

**Hidden Sun  $i$  of row  $l$  ( $l = 2, 3, \dots, K-1$ ):**

Inputs Used:  $w_{li0} z_{(l-1)0}, w_{li1} z_{(l-1)1}, \dots, w_{liM_{(l-1)}} z_{(l-1)M_{(l-1)}}$   
Local Memory Value Used: None  
Output:  $z_{li} = s(\sum_j w_{lij} z_{(l-1)j})$   
Local Memory Value Stored:  $I_{li} \equiv \sum_j w_{lij} z_{(l-1)j}$

**Output Sun  $i$  of row  $l = K$ :**

Inputs Used:  $w_{Ki0} z_{(K-1)0}, w_{Ki1} z_{(K-1)1}, \dots, w_{KiM_{(K-1)}} z_{(K-1)M_{(K-1)}}$   
Local Memory Value Used: None  
Output:  $z_{Ki} = y_i' = s(\sum_j w_{Kij} z_{(K-1)j})$   
or, optionally,  $z_{Ki} = y_i = \sum_j w_{Kij} z_{(K-1)j}$   
Local Memory Values Stored:  $I_{Ki} \equiv \sum_j w_{Kij} z_{(K-1)j}, y_i'$

where  $s(I) = 1/(1 + e^{-I})$  is the *sigmoid function* of the network,  $s' = ds(I)/dI$  is the first derivative of the sigmoid function,  $M_l$  is the number of units on row  $l$ ,  $w_{lij}$  is the weight of planet  $j$  of unit  $i$  of row  $l$ ,  $z_{li} = z_i$  (where  $z_i$  is the  $i^{\text{th}}$  component of the input vector  $\mathbf{x}$ ), and the network's output signals  $z_{Ki}$  are equal to the components  $y_i'$  of the network's output vector  $\mathbf{y}'$  – the network's estimate of the “correct” or “desired” output  $\mathbf{y}$  (the components of which are supplied to the network for use in learning). Note that the output signals of the output suns can either be sigmoided or unsigmoided, as desired by the user. Further note that primes (') are used to indicate both the correct output signals of the network as well as the first derivative of the sigmoid function.

Table 1: Processing element transfer functions for the backpropagation neural network architecture. Three types of processing element are used (planets, hidden suns, and output suns). For each processing element type the transfer function used on the forward pass of the network and the backward pass of the network is given.

This process is continued until the network reaches a satisfactory level of performance (which is usually tested using data that was not used during training, but which was drawn at random in the same way as the training data); or until the user gives up and decides to try a new starting weight set or a new network configuration (see Section 5 below for a discussion of backpropagation network training). After successfully training the network to a suitably low level of error the network can then be put through further operational testing to qualify it for deployment. Since training is typically not used in the deployed version, the backward pass can be eliminated – allowing a considerable computational savings to be realized.

### 3 Backpropagation Error Surfaces

Given a function  $f$ , an associated  $\mathbf{x}$ -selection probability density function  $\rho$ , and an associated backpropagation architecture (i.e., with a given number of layers  $K$ , and a given number of units  $M_l$  per hidden layer, with either linear or sigmoided output units) intended to approximate  $f$ , we now define a means for measuring the accuracy of this approximation as a function of the network's weights and discuss the resulting *error surface*.

First, define  $\mathbf{w}$  (the *weight vector* of the network) to be the vector with components consisting of the weights of all of the planets of the network, starting with the weight of the first planet of the first processing element of layer 2 (the first hidden layer) and ending with the weight of the last planet of the last processing element of the output layer  $K$ . To make things simple we shall refer to the components of  $\mathbf{w}$  as  $w_1, w_2, \dots$ ; rather than  $w_{210}, w_{211}, \dots$ . To make our notation simple the vector  $(z_{K1}, z_{K2}, \dots, z_{Km})$  (the network's estimate  $\mathbf{y}'$  of the correct output  $\mathbf{y}$ ) shall be called  $\mathbf{B}$ . Note that  $\mathbf{B}$  is a function of the input vector  $\mathbf{x}$  and the network weight vector  $\mathbf{w}$ . Thus, we will write  $\mathbf{B}(\mathbf{x}, \mathbf{w})$ . Let  $(\mathbf{x}_k, \mathbf{y}_k)$  be the example used on the  $k^{\text{th}}$  testing trial (i.e.,  $\mathbf{y}_k = f(\mathbf{x}_k)$ ). As before, the  $\mathbf{x}_k$ 's are drawn from  $A$  in accordance with a fixed probability density function  $\rho$ .

Given the above, let  $F_k = |f(\mathbf{x}_k) - \mathbf{B}(\mathbf{x}_k, \mathbf{w})|^2$ .  $F_k$  is the square of the approximation error made on the  $k^{\text{th}}$  testing trial. For the purposes of the discussion below we shall assume that  $\mathbf{w}$  is fixed (i.e., batch\_size is set to  $\infty$  to shut off learning). Then we define  $F(\mathbf{w})$  to be

$$F(\mathbf{w}) \equiv \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N F_k. \quad (1)$$

Because the mapping  $f$  is bounded the set of all  $\mathbf{y}_k$ 's are bounded. Thus, it is easy to show that this limit exists because (for a fixed  $\mathbf{w}$ )  $\mathbf{B}$  is a continuous mapping from the compact set  $A$  into  $\mathbf{R}^m$ , and thus the set of all  $\mathbf{B}$ 's are bounded also. Thus, the variance of the random variable  $F_k$  is bounded. So, by Kolmogorov's probability theorem [4], the random variable  $F_k$  obeys the strong law of large numbers. Therefore, the above  $F$  sum must almost surely converge to the expected value  $F(\mathbf{w})$  of  $F_k$ . We call  $F(\mathbf{w})$  the *mean squared error function* of the network; which is often just shortened to *error function*. Note that  $F(\mathbf{w}) \geq 0$  because  $F$  is the average of non-negative quantities.

#### Backward Pass

Input Used:  $z_{li} (= \delta_{li})$   
Weight Value Used:  $w_{lij}$   
Local Memory Value Used:  $\Delta_{li}$   
Output:  $w_{lij}^{\text{old}} \delta_{li}$   
Weight and Local Memory Value Update:  
IF (count = batch\_size)

THEN {  
 $w_{lij}^{\text{new}} = w_{lij}^{\text{old}} + \frac{\alpha \Delta_{li}}{\text{batch\_size}}$   
 $\Delta_{lij}^{\text{new}} = \delta_{li} z_{(l-1)j}$   
count = 1 }

ELSE {  
 $w_{lij}^{\text{new}} = w_{lij}^{\text{old}}$   
 $\Delta_{lij}^{\text{new}} = \Delta_{lij}^{\text{old}} + \delta_{li} z_{(l-1)j}$   
count = count + 1 }

Inputs Used:  $w_{(l+1)1} \delta_{(l+1)1}, w_{(l+1)2} \delta_{(l+1)2}, \dots, w_{(l+1)M_l} \delta_{(l+1)M_l}$   
Local Memory Value Used:  $I_{li}$   
Output:  $z_{li} = \delta_{li}$   
 $\equiv s'(I_{li}) \sum_l w_{(l+1)li} \delta_{(l+1)l}$   
Local Memory Value Stored: None

Input Used:  $y_i$   
Local Memory Value Used:  $I_{Ki}$   
Output:  $z_{Ki} = \delta_{Ki}$   
 $\equiv s'(I_{Ki}) (y_i - y_i')$   
Local Memory Value Stored: None

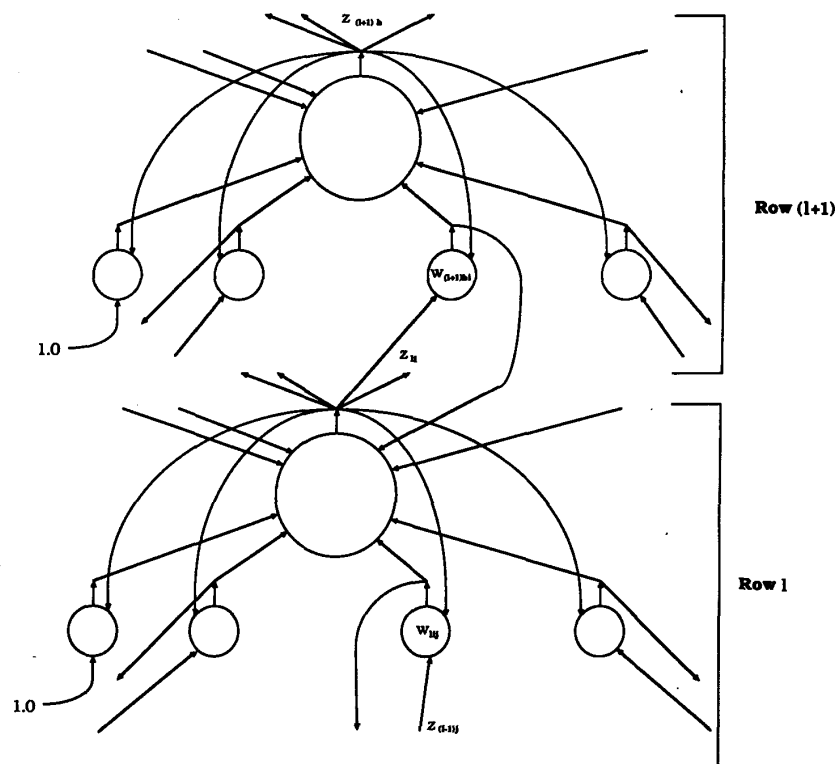


Figure 2: Architectural detail and interaction of two backpropagation network units on two adjacent rows (row  $l$  below and row  $l+1$  above). Each unit consists of a *sun* and several *planets*. Each hidden row sun sends its output to one planet of each unit of the next row. Each planet that receives an input from the previous row (the *bias* planets all receive the constant input  $z_{(l-1)0} \equiv 1.0$  from a single processing element which is not shown) sends a connection back to the sun that supplied that input. It is through these paths that the products of the errors times the weights ( $w_{ij}\delta_i$ ) are back-propagated.

The *error surface* of a backpropagation network is the surface defined by the equation  $F = F(\mathbf{w})$  in the  $Q+1$ -dimensional space of vectors  $(\mathbf{w}, F)$ , where  $Q$  is the number of dimensions in the vector  $\mathbf{w}$  (i.e., the number of planets in the network). The variable  $\mathbf{w}$  ranges over its  $Q$ -dimensional space and for each  $\mathbf{w}$  a non-negative surface height  $F$  is defined by  $F(\mathbf{w})$ . In other words, given any selection of weights  $\mathbf{w}$ , the network will make an average squared error  $F(\mathbf{w})$  in its approximation of the function  $f$ . We now consider the shape of this error surface.

As will be shown below, the generalized delta rule learning law used with the backpropagation neural network has the property that, given any starting point  $\mathbf{w}_0$  on the error surface that is not a minimum, the learning law will modify the weight vector  $\mathbf{w}$  so that  $F(\mathbf{w})$  will decrease. In other words, the learning law uses examples provided during training to decide how to modify the weight vector so that the network will do a better job of approximating the function  $f$ . Given this behavior, the next issue is to assess how valuable this property is.

Until recently, the shape of backpropagation error surfaces was largely a mystery. Two basic facts have emerged so far. First, experience has shown that many backpropagation error surfaces are dominated by flat areas and troughs that have very little slope. In these areas it is necessary to move the weight value a considerable distance before a significant drop in error occurs. Since the slope is shallow it turns out that the generalized delta rule has a hard time determining which way to move the weight to reduce the error. Often, great numerical precision (e.g., 32-bit floating point) and patience must be employed to make significant progress.

The other basic fact about error surfaces that has emerged concerns the existence of local minima. Until recently, it was not known for certain whether backpropagation error surfaces have local minima at error levels above the levels of the global minima of the surfaces (due to weight permutations there are always many global minima). Experience suggested that such minima might not exist because usually when training failed to make downhill progress and the error level was high it was discovered that further patience (or the use of one of the training augmentations described in Section 5 below) would eventually lead to the weight moving away from what was clearly a shallow spot on the surface and onto a steeper part. Thus, it was somewhat surprising when McInerney, Haines, Biafore, and Hecht-Nielsen [36] discovered a local minimum (at a very high error level) in a backpropagation error surface in June 1988. Finding this local minimum was not easy. It required the use of the Symbolics Corporation's Macsyma<sup>TM</sup> mathematics expert system running on the largest Digital Equipment Corporation VAX<sup>TM</sup> computer to generate approximately a megabyte of FORTRAN code (the closed-form formulas for the error function and all of its first and second partial derivatives with respect to the weights) and a 12-hour run on a Cray-2 supercomputer (of a program that called this error surface code) to accomplish this discovery. The existence of the minimum was proven by showing that all of the first partial derivatives of the mean squared error function went to zero at this point and that the Hessian of this function (the matrix of second partial derivatives) was strongly positive-definite at this point.

In summary, three basic facts are known about backpropagation error surfaces. First, because of combinatoric permutations of the weights that leave the network input-output function unchanged, these functions typically have huge numbers of global minima (which may lie at infinity for some problems). This causes the error surfaces to be highly degenerate and to have numerous 'troughs'. Secondly, error surfaces have a multitude of areas with shallow slopes in multiple dimensions simultaneously. These typically occur because particular combinations of weights cause the weighted sums of one or more suns (with sigmoidal outputs) to be large in magnitude. When this occurs the output of that sun (and therefore the value of  $F$ ) is insensitive to small weight changes, since these simply move the weighted sum value back and forth along one of the shallow tails of the sigmoid function. Thirdly, it is now established that local minima do actually exist. Beyond these three facts, little is known. How many non-global minima are there, compared to the number of global minima? Are local minima excluded from regions near global minima? How large are the attractive basins of each of the different types of minima? What forms do the boundaries between attractive basins take? Can some of the nonlinear optimisation and unconstrained optimization techniques developed by researchers in operations research and mathematical economics be usefully applied? Clearly, more research into error surfaces is needed.

## 4 Function Approximation with Backpropagation

The question of what functional forms can be approximated by neural networks has had pivotal importance in the history of neurocomputing. For example, the fact that a narrowly formulated type of perceptron could be shown incapable of implementing the EXCLUSIVE OR logical operation [38] was used in the 1960's as an argument to divert funding from neurocomputing to artificial intelligence. Recently, similar questions have been raised and claims that little progress has been made on this front [38] have generated concern.

Clear insight into the versatility of neural networks for use in function approximation came with the discovery [20] that a classic mathematical result of Kolmogorov [30] was actually a statement that for any continuous mapping  $f: [0, 1]^n \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  there must exist a three-layer neural network (having an input or "fanout" layer with  $n$  processing elements, a hidden layer with  $(2n+1)$  processing elements, and an output layer with  $m$  processing elements) that implements  $f$  exactly. This result gave hope that neural networks would turn out to be able to approximate any function that arises in the real world.

Kolmogorov's theorem was a first step. The following new result shows that the backpropagation network is itself able to implement any function of practical interest to any desired degree of accuracy. To make the exposition precise, we start with some background information.

Let  $[0, 1]^n$  be the closed unit cube in  $n$ -dimensional Euclidean space. Given any square-integrable function  $g: [0, 1]^n \subset \mathbb{R}^n \rightarrow \mathbb{R}$  (i.e.,  $\int_{[0, 1]^n} |g(x)|^2 dx$  exists), it can be shown by the theory of Fourier series [12] that the series

$$F(g, x, N) = \sum_{k_1=-N}^N \sum_{k_2=-N}^N \cdots \sum_{k_n=-N}^N c_{k_1, k_2, \dots, k_n} \exp \left( 2\pi i \sum_{q=1}^n k_q x_q \right) = \sum_{\mathbf{k}} c_{\mathbf{k}} \exp(2\pi i \mathbf{k} \cdot \mathbf{x}) \quad (2)$$

where

$$c_{k_1, k_2, \dots, k_n} = c_{\mathbf{k}} = \int_{[0, 1]^n} g(\mathbf{x}) \exp(-2\pi i \mathbf{k} \cdot \mathbf{x}) d\mathbf{x} \quad (3)$$

converges to  $f$  in the sense that

$$\lim_{N \rightarrow \infty} \int_{[0, 1]^n} |g(\mathbf{x}) - F(g, \mathbf{x}, N)|^2 d\mathbf{x} = 0. \quad (4)$$

This is an example of the property of having the integral of the square of the error of approximation of one function by another go to zero. This property is described by the statement that the approximation can be achieved to within any desired degree of accuracy in the *mean squared error sense*. This leads to a definition.

Given a function  $f: [0, 1]^n \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , we say that  $f$  belongs to  $L_2$  (or "is  $L_2$ ") if each of  $f$ 's coordinate functions is square-integrable on the unit cube. For functions of this class it is assumed that the  $\mathbf{x}$  vectors are chosen uniformly in  $[0, 1]^n$  (relaxing this condition is fairly easy). Clearly if a vector function of a vector variable is in  $L_2$  then each of its components can be approximated by its Fourier series to any desired degree of accuracy in the mean squared error sense. With this background as preamble, the new result is now presented.

**Theorem 1** Given any  $\epsilon > 0$  and any  $L_2$  function  $f: [0, 1]^n \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , there exists a three-layer backpropagation neural network that can approximate  $f$  to within  $\epsilon$  mean squared error accuracy.

**Proof:** Let  $\epsilon$  be the accuracy to which we wish to approximate each of the coordinate functions  $f_i$  (where  $f(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))$ ) of  $f$  by using a three-layer backpropagation neural network. Note that, by virtue of the results from Fourier series theory cited above, given any  $\delta_1 > 0$ , there exists a positive integer  $N$  and coefficients  $c_{\mathbf{k}}$  such that

$$\int_{[0, 1]^n} |f_i(\mathbf{x}) - \sum_{\mathbf{k}} c_{\mathbf{k}} \exp(2\pi i \mathbf{k} \cdot \mathbf{x})|^2 d\mathbf{x} < \delta_1. \quad (5)$$

We begin by showing that each of the sine and cosine terms required in the Fourier series (i.e., in the real part of the complex Fourier series - the imaginary part vanishes) for  $f_i$  can be implemented to any desired degree of absolute accuracy by a subset of a three-layer backpropagation neural network. The idea is to use the input layer units, a contiguous subset of the hidden layer units, and the corresponding portion of a single output layer unit - unit  $l$  out of a total of  $m$  output units - to implement each sine or cosine function. The output units of the backpropagation neural network used here do not employ sigmoid functions on their outputs. To carry out this approximation, first note that the input layer, a subset  $H$  of the hidden layer (assumed here to be comprised of contiguous hidden units), and the  $l^{\text{th}}$  output unit can compute any sum of the form

$$\sum_{i \in H} v_{li} s \left( \sum_{j=0}^n w_{ij} x_j \right) \quad (6)$$

where  $w_{0j} \equiv 0$  for  $j = 1, 2, \dots, n$ , and  $x_0 \equiv 1$  is the bias input to each hidden unit. Next, we note that each of the arguments of the sine and cosine functions in the Fourier approximation of  $f_i$  (namely the terms  $u = u(\mathbf{k}, \mathbf{x}) = 2\pi i \mathbf{k} \cdot \mathbf{x}$ ) are of the form

$$\sum_{j=0}^n w_{ij} x_j. \quad (7)$$

Since, by simply adjusting the argument bias weight  $w_{i0}$  by  $-(\pi/2)$  we can change a sine into a cosine, we shall concern ourselves exclusively with sines. Thus, to absolutely approximate a particular Fourier series sine or cosine function all we need to show is that, given any  $\delta_2 > 0$ , we can find coefficients  $v_{li}$  and  $w_{ij}$  such that

$$|\sin(u(\mathbf{k}, \mathbf{x})) - \sum_{i \in H} v_{li} s\left(\sum_{j=0}^n w_{ij} x_j\right)| < \delta_2 \quad (8)$$

for each  $\mathbf{x} \in [0, 1]^n$ . To show this, choose the  $w_{ij}$  such that

$$\sum_{j=0}^n w_{ij} x_j = \beta_i (u(\mathbf{k}, \mathbf{x}) - \alpha_i) \quad (9)$$

where the quantities  $\beta_i$  and  $\alpha_i$  are arbitrary real constants to be selected below. We are going to dedicate a certain (typically large) subset  $H_{\mathbf{k}}$  of hidden layer units (assumed to be contiguous) to the calculation of each  $\sin(u(\mathbf{k}, \mathbf{x}))$  and  $\cos(u(\mathbf{k}, \mathbf{x}))$ . So, with this transformation for each hidden layer unit within  $H_{\mathbf{k}}$  we can rewrite the above inequality as

$$|\sin(u(\mathbf{k}, \mathbf{x})) - \sum_{i \in H_{\mathbf{k}}} v_{li} s(\beta_i (u(\mathbf{k}, \mathbf{x}) - \alpha_i))| < \delta_2 \quad (10)$$

It is easy to show, given a sufficient number of hidden layer units in  $H_{\mathbf{k}}$ , that this inequality can always be satisfied. The first step in this demonstration is to examine the geometrical form of the sum

$$S(\alpha, \beta, \mathbf{v}_l, \mathbf{x}) \equiv \sum_{i \in H_{\mathbf{k}}} v_{li} s(\beta_i (u - \alpha_i)) \quad (11)$$

which, by the above argument, can be computed by the input layer, hidden layer subset  $H_{\mathbf{k}}$ , and output unit  $l$ . Note that for  $\mathbf{x} \in [0, 1]^n$ ,  $u$  ranges over some closed interval  $d(\mathbf{k}) \leq u \leq e(\mathbf{k})$ . It is a closed interval because the affine transformation  $\mathbf{x} \mapsto u$  is continuous and continuous maps preserve compactness and simple connectedness. So, we need to approximate the function  $\sin(u)$  on the closed interval  $d \leq u \leq e$  using this sum. To do this, partition the interval as follows

$$d = \alpha_{i_1} < \alpha_{i_2} < \alpha_{i_3} < \dots < \alpha_{i_{p+1}} = e \quad (12)$$

where  $i_{p+1} = i_p + 1$  and  $\cup_{p=1}^{p+1} \{i_p\} = H_{\mathbf{k}}$ .

By letting  $\beta_{i_1} = 0$  and by choosing the  $\beta_{i_p}$ ,  $p > 1$  to be sufficiently large it is easy to see that the sum  $S(\alpha, \beta, \mathbf{v}_l, \mathbf{x})$  has the geometrical form shown in Figure 3, since each of the terms in the sum is a steep sigmoid (sort of like a step function with microscopically rounded corners) that is essentially equal to 1 for  $u > \alpha_{i_p}$ , and equal to 0 for  $u < \alpha_{i_p}$ , (and equal to 0.5 when  $u = \alpha_{i_p}$ ).

Clearly,  $S$  has the approximate form of a staircase, where the step lengths are determined by the differences  $(\alpha_{i_{p+1}} - \alpha_{i_p})$  and the step heights are determined by the coefficients  $v_{li}$ . By setting the sigmoid gains  $\beta_{i_p}$ ,  $p > 1$  to high enough values this basic staircase form can always be achieved, no matter how small the steps become or how many steps are used.

Given the above facts about the geometrical form of the sum  $S$ , Figure 4 demonstrates graphically that no matter how small  $\delta_2 > 0$  is chosen, a sum  $S(\alpha, \beta, \mathbf{v}_l, \mathbf{x})$  can be constructed so that it always remains within the  $\delta_2$  error band around  $\sin(u)$ . By starting at the left end of the interval  $[d, e]$ , and working to the right, successive sigmoid stairsteps can be incrementally added to achieve this.

Thus, we have shown that by choosing  $H_{\mathbf{k}}$  to have an adequately large number of processing elements and by properly selecting the  $\alpha$ ,  $\beta$ , and  $\mathbf{v}_l$  vectors, the partial sum output (due to the members of  $H_{\mathbf{k}}$ ) will be within  $\delta_2$  of  $\sin(u(\mathbf{k}, \mathbf{x}))$  (or  $\cos(u(\mathbf{k}, \mathbf{x}))$ , if that is what is being fitted), for all  $\mathbf{x} \in [0, 1]^n$ . Since output unit  $l$  receives the inputs from the hidden layer units of each of the sine and cosine subsets  $H_{\mathbf{k}}$  for every  $\mathbf{k} \in V \equiv \{-N, -N+1, \dots, N\}^n$  we can then approximately generate the Fourier series for  $f_l$  by simply multiplying all of the sine (or cosine) coefficients  $v_{li}$  for  $i \in H_{\mathbf{k}}$  by the appropriate combinations of the real and imaginary parts of  $c_{l\mathbf{k}}$  (the complex Fourier expansion coefficient of  $f_l$  for term  $\mathbf{k}$  in the expansion). Call these sine and cosine multipliers  $a(l, \mathbf{k})$  and  $b(l, \mathbf{k})$ , respectively, and let  $y_l(\mathbf{x})$  be the output signal of output unit  $l$ . Thus, we get

$$F(f_l, \mathbf{x}, N) - y_l(\mathbf{x}) = \sum_{\mathbf{k} \in V} a(l, \mathbf{k}) [\sin(2\pi i \mathbf{k} \cdot \mathbf{x}) - S(\alpha, \beta, \mathbf{v}_l, \mathbf{x})] + b(l, \mathbf{k}) [\cos(2\pi i \mathbf{k} \cdot \mathbf{x}) - S'(\alpha, \beta, \mathbf{v}_l, \mathbf{x})], \quad (13)$$

where the  $S'$  sum is that used for the  $H_{\mathbf{k}}$  cosine term.

Putting all of the above together we get:

$$\begin{aligned} \int_{[0,1]^n} |f_l(\mathbf{x}) - y_l(\mathbf{x})|^2 d\mathbf{x} &= \int_{[0,1]^n} |f_l(\mathbf{x}) - F(f_l, \mathbf{x}, N) + F(f_l, \mathbf{x}, N) - y_l(\mathbf{x})|^2 d\mathbf{x} \\ &\leq \int_{[0,1]^n} |f_l(\mathbf{x}) - F(f_l, \mathbf{x}, N)|^2 d\mathbf{x} + \int_{[0,1]^n} |F(f_l, \mathbf{x}, N) - y_l(\mathbf{x})|^2 d\mathbf{x} \\ &< \delta_1 + (\delta_2)^2 \sum_{\mathbf{k} \in V} ([a(l, \mathbf{k})]^2 + [b(l, \mathbf{k})]^2). \end{aligned} \quad (14)$$

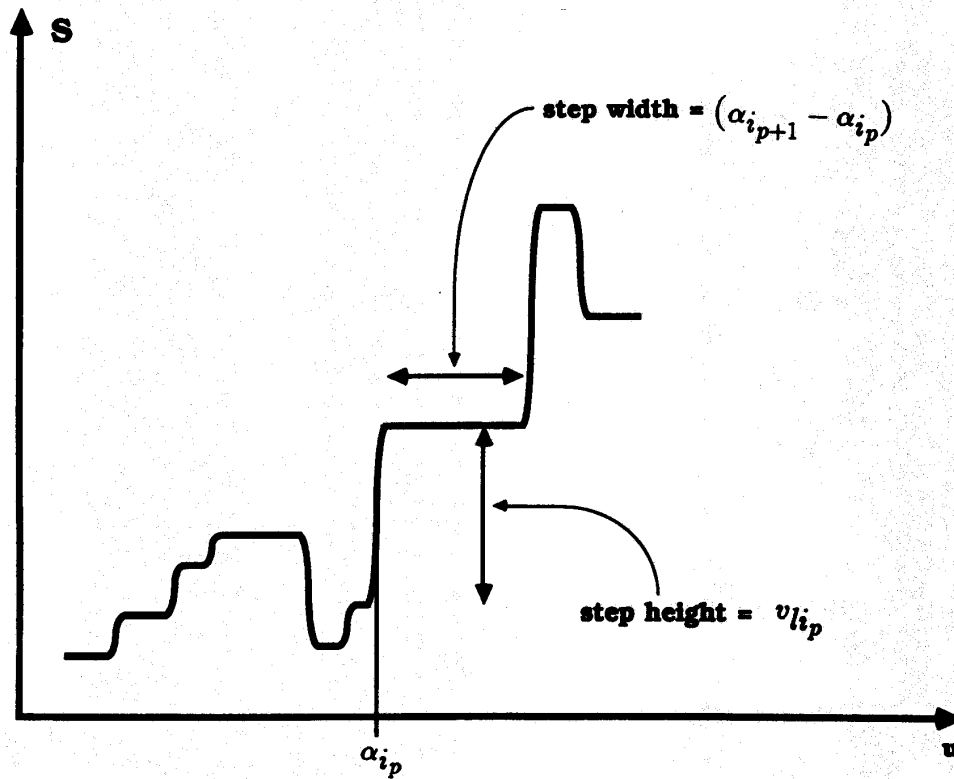


Figure 3: Geometrical Form of the function  $S(\alpha, \beta, v_l, x)$  with  $\beta_l = 0$  and large  $\beta_p$ , for  $p > 1$ .

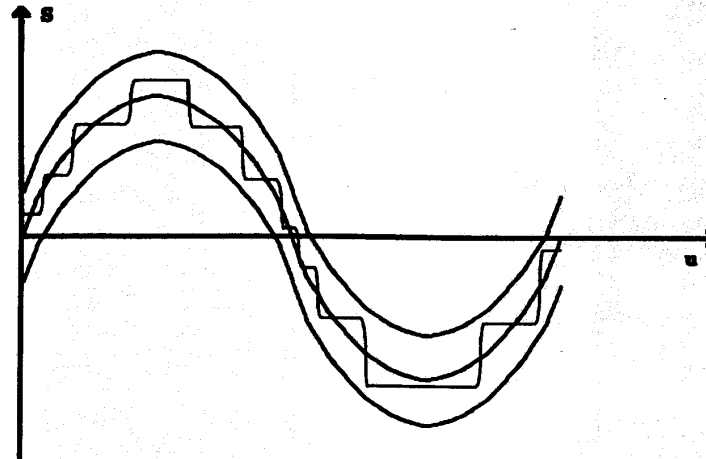


Figure 4: Graphical demonstration that  $\alpha$ ,  $\beta$ , and  $v_l$  vectors can always be chosen so that the function  $S(\alpha, \beta, v_l, x)$  will always remain within an error band of width  $\delta_2$  around  $\sin(u)$ .

The first inequality follows from Pythagorus' theorem and the second from the facts that  $f_l$  can be approximated to within  $\delta_1$  by its Fourier series and that the absolute differences between  $\sin(2\pi i k \cdot x)$  and  $S(\alpha, \beta, v_l, x)$  are always less than  $\delta_2$ , for all  $x \in [0, 1]^n$ , followed by a second application of Pythagorus' theorem. So, given any desired mean squared error  $\epsilon > 0$ , we can select  $\delta_1$  and  $\delta_2$  sufficiently small so that we will have

$$\int_{[0,1]^n} |f_l(x) - y_l'(x)|^2 dx < \frac{\epsilon}{n}, \quad (15)$$

or

$$\int_{[0,1]^n} |f(\mathbf{x}) - y'(\mathbf{x})|^2 d\mathbf{x} < \epsilon \quad (16)$$

which proves the theorem. To relate this result to the error function of Section 3, we need only note that

$$F(\mathbf{w}) = \int_{[0,1]^n} |f(\mathbf{x}) - y'(\mathbf{x})|^2 d\mathbf{x} \quad (17)$$

for  $L_2$  functions.

The space  $L_2$  includes every function that could ever arise in a practical problem. For example, it includes the continuous functions and it includes all discontinuous functions that are piecewise continuous on a finite number of subsets of  $[0, 1]^n$ .  $L_2$  also contains much nastier functions than these, but they are only of mathematical interest.

It is important to realize that although this theorem proves that 'three layers are always enough', in solving real-world problems it is often essential to have four, five, or even more layers. This is because for many problems an approximation with three layers would require an impractically large number of hidden units, whereas an adequate solution can be obtained with a tractable network size by using more than three layers. Thus, the above result should not influence the process of selecting an appropriate backpropagation architecture for a real-world problem - except to provide the confidence that comes from knowing that an appropriate backpropagation architecture *must* exist.

Finally, although the above theorem guarantees the ability of a multilayer network *with the correct weights* to accurately implement an arbitrary  $L_2$  function, it does not comment on whether or not these weights can be *learned* using any existing learning law. That is an open question.

## 5 Backpropagation Learning Laws

Most, if not all, of the widely disseminated published derivations of the backpropagation neural network architecture learning law - the generalized delta rule - are flawed or incomplete. (Chapter 8 of [45] provides one of the better derivations, although a highly idiosyncratic notation is used.) This Section provides a careful derivation of the generalized delta rule and a brief discussion of some current research directions in backpropagation learning laws.

The first step in the derivation of the generalized delta rule is to recall that our goal is to move the  $\mathbf{w}$  vector in a direction so that the value of  $F$  will be smaller at the new value. Obviously, the best approach would be to move in the  $\mathbf{w}$  direction in which  $F$  is decreasing most rapidly. Assuming that  $F$  is differentiable, ordinary vector analysis tells us that this direction of maximum decrease is given by  $-\nabla_{\mathbf{w}} F(\mathbf{w}) = -\left(\frac{\partial F}{\partial w_1}, \frac{\partial F}{\partial w_2}, \dots, \frac{\partial F}{\partial w_Q}\right)$ , where, as in Section 3, it is assumed that there are  $Q$  weights in the network (i.e., the vector  $\mathbf{w}$  is assumed to have  $Q$  components). However, before proceeding with this line of thought it is essential to show that  $F$  is in fact differentiable.

First, note that because backpropagation networks are made up entirely of affine transformations and smooth sigmoid functions,  $B(\mathbf{x}, \mathbf{w})$  is a  $C^\infty$  function of both  $\mathbf{x}$  and  $\mathbf{w}$  (i.e., its partial derivatives of all orders with respect to the components of  $\mathbf{x}$  and  $\mathbf{w}$  exist and are continuous) and its derivative limits all converge uniformly on compact sets. Thus,  $F_k$ , as a function of  $\mathbf{w}$ , inherits these same properties, since  $F_k = |f(\mathbf{x}_k) - B(\mathbf{x}_k, \mathbf{w})|^2$ , and since  $f(\mathbf{x}_k)$  is a constant vector. What we will need later on are two facts: first, that  $F$  is a differentiable function, and second, that  $\nabla_{\mathbf{w}} F(\mathbf{w}) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \nabla_{\mathbf{w}} F_k(\mathbf{w})$  almost surely. These two facts can be simultaneously established. First, fix  $\mathbf{w}$  and, given an arbitrary example sequence  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots\}$  (as defined in Section 3 above), let

$$G_i(\Gamma, \Delta) = \frac{1}{[1/\Gamma]} \sum_{k=1}^{[1/\Gamma]} \left( \frac{F_k(\mathbf{w} + \Delta \sigma_i) - F_k(\mathbf{w})}{\Delta} \right), \quad (18)$$

where  $\sigma_i$  is the unit basis vector along the  $i^{\text{th}}$  coordinate axis in weight space, where the real variables  $\Gamma$  and  $\Delta$  both range over a compact neighborhood  $U$  of zero (with zero removed) and where  $[u]$  (the floor function) is the largest integer less than or equal to  $u$ . Because the derivatives of  $F_k$  all converge uniformly on compact sets in weight space the limit

$$\lim_{\Delta \rightarrow 0} G_i(\Gamma, \Delta) = [1/\Gamma] \sum_{k=1}^{[1/\Gamma]} \frac{\partial F_k}{\partial w_i} \quad (19)$$

converges uniformly on  $U$ . Further, as shown in Section 2 above, for each  $\Delta$  in  $U$  the limit

$$\lim_{\Gamma \rightarrow 0} G_i(\Gamma, \Delta) = \left( \frac{F(\mathbf{w} + \Delta \sigma_i) - F(\mathbf{w})}{\Delta} \right), \quad (20)$$

converges almost surely. Thus, by the elementary theory of iterated limits [31] the double limits

$$\lim_{\Gamma \rightarrow 0} \lim_{\Delta \rightarrow 0} G_i(\Gamma, \Delta) \quad (21)$$

and

$$\lim_{\Delta \rightarrow 0} \lim_{\Gamma \rightarrow 0} G_i(\Gamma, \Delta) \quad (22)$$

both exist and are equal, almost surely. Thus, we have shown that, almost surely,  $F$  is a differentiable function and that  $\nabla_{\mathbf{w}} F(\mathbf{w}) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \nabla_{\mathbf{w}} F_k(\mathbf{w})$ .

Given these preliminaries, the generalized delta rule is now derived. The basic idea is to derive a formula for calculating  $-\nabla_{\mathbf{w}} F(\mathbf{w}) = -\left(\frac{\partial F}{\partial w_1}, \frac{\partial F}{\partial w_2}, \dots, \frac{\partial F}{\partial w_Q}\right)$  by using training examples. To do this, we focus on the calculation of  $\frac{\partial F}{\partial w_i}$ . We first note that, using the above results,



$$\frac{\partial F(\mathbf{w})}{\partial w_p} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \frac{\partial F_k(\mathbf{w})}{\partial w_p} \quad (23)$$

but, reverting to full weight indices (where the  $p^{\text{th}}$  weight has full index  $lij$ ),

$$\frac{\partial F_k(\mathbf{w})}{\partial w_p} = \frac{\partial F_k(\mathbf{w})}{\partial w_{lij}} = \frac{\partial F_k}{\partial I_{li}} \frac{\partial I_{li}}{\partial w_{lij}}, \quad (24)$$

because any functional dependence of  $F_k$  on  $w_{lij}$  must be through  $I_{li}$  (see Section 2). If we now define  $\delta_{li}^k$  to be  $\frac{\partial F_k}{\partial I_{li}}$  and evaluate  $\frac{\partial I_{li}}{\partial w_{lij}}$  in terms of the formula for  $I_{li}$ , we get

$$\frac{\partial F_k(\mathbf{w})}{\partial w_p} = \delta_{li}^k \frac{\partial}{\partial w_{lij}} \left( \sum_q w_{liq} z_{(l-1)q}^k \right) = \delta_{li}^k z_{(l-1)j}^k, \quad (25)$$

where  $z_{lj}^k$  is the output signal of the  $j^{\text{th}}$  sun of the  $l^{\text{th}}$  row on the forward pass of the  $k^{\text{th}}$  training trial. But,

$$\delta_{li}^k \equiv \frac{\partial F_k}{\partial I_{li}} = \frac{\partial F_k}{\partial z_{li}} \frac{\partial z_{li}}{\partial I_{li}} = \frac{\partial F_k}{\partial z_{li}} s'(I_{li}) \quad (26)$$

where  $s'()$  is the ordinary derivative of the sigmoid function  $s$  with respect to its argument. If row  $l$  is the output layer, then

$$\frac{\partial F_k}{\partial z_{li}} = \frac{\partial}{\partial z_{li}} \sum_{p=1}^m (y_p^k - z_{lp}^k)^2 = -2 (y_i^k - z_{li}^k), \quad (27)$$

where  $y_p^k$  is the  $p^{\text{th}}$  component of  $\mathbf{y}_k$ . Thus, for the output layer units

$$\delta_{li}^k = -2 (y_i^k - z_{li}^k) s'(I_{li}). \quad (28)$$

If row  $l$  is a hidden layer, then

$$\frac{\partial F_k}{\partial z_{li}} = \sum_{p=1}^{M_{(l+1)}} \frac{\partial F_k}{\partial I_{(l+1)p}} \frac{\partial I_{(l+1)p}}{\partial z_{li}}. \quad (29)$$

Thus, for the hidden layer units

$$\delta_{li}^k = s'(I_{li}) \sum_p \delta_{(l+1)p}^k w_{(l+1)pi}. \quad (30)$$

So, substituting these results back into the formula for  $\frac{\partial F}{\partial w_{lij}}$ , gives

$$\frac{\partial F}{\partial w_p} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \frac{\partial F_k}{\partial w_{lij}} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \delta_{li}^k z_{(l-1)j}^k. \quad (31)$$

Thus, to move  $\mathbf{w}$  in a direction that will decrease  $F$  (and thus improve the accuracy by which the network approximates the function  $f$ ), all we need do is move in the direction  $-\nabla_{\mathbf{w}} F(\mathbf{w})$ . The amount to be moved must be small, to avoid the problem of overshooting the bottom of the hill. Thus, a reasonable learning law is

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \alpha \nabla_{\mathbf{w}} F(\mathbf{w}), \quad (32)$$

or

$$w_{lij}^{\text{new}} = w_{lij}^{\text{old}} - \alpha \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \delta_{li}^k z_{(l-1)j}^k \quad (33)$$

where  $\alpha > 0$  is a small constant called the *learning rate*. This is the *generalized delta rule* learning law.

The rigorous generalized delta rule is approximated by the learning law used in the backpropagation architecture presented in Section 2, which substitutes a finite averaging sum of *batch\_size* terms for the infinite-size average of the formula. Since we know that the limit converges this will clearly be acceptable, assuming *batch\_size* is chosen large enough.

A number of variants of this learning law (and of other parts of the backpropagation neural network architecture) have been presented [9,25]. One learning law variant in common use is

$$w_{lij}^{\text{new}} = w_{lij}^{\text{old}} - \alpha \delta_{li}^k z_{(l-1)j}^k. \quad (34)$$

This law is the backpropagation analog of the *Widrow-Hoff LMS rule* [54,55,56]. The proof that this law also carries out gradient descent on the backpropagation error surface was provided by Hirsch [22]. Another common variant is the *momentum* version of the above law (see [45] for details).

All of the above learning laws suffer from the same basic problem: they move downhill in short jumps. Thus, even if the network's initial weight vector lies within the attractive basin of a global minimum of the error surface, getting to the bottom can take many jumps.

The development of new learning laws for backpropagation continues. The general goal is to provide a faster descent to the bottom of the error surface. One line of investigation [48,42,3] is exploring the use of approximations to the pseudoinverse [29] of the Hessian matrix ( $\partial^2 F / \partial w_i \partial w_j$ ) to calculate individually variable learning rates ( $\alpha$  values in Equation (34)) for each weight  $w_{ij}$ . The underlying concept is to use Newton's method for finding a place where  $\nabla_{\mathbf{w}} F(\mathbf{w}) = 0$ . So far, this work shows promise, but a major advance in convergence speed has yet to be realized for arbitrary problems. Another approach that shows promise in early tests is the incremental network growing technique of Ash [2]. This method starts with a small number of units in each hidden layer and adds additional units in response to changes in the error level (as measured using test data).

It is probably reasonable to anticipate that faster learning techniques will be developed in the future as more becomes known about the structure of backpropagation error surfaces, and as these facts are exploited. The advent of such faster learning techniques will correspondingly increase the breadth of applicability of the backpropagation architecture.

Finally, it is important to note that this paper has examined only one (albeit perhaps the most important) variant of the backpropagation neural network architecture and its learning law. Other important variants exist. Examples of these include architectures in which connections skip layers [45], recurrent architectures [45], "sigma-pi" higher order architectures [45], the "second-order" learning laws of Parker [39], the "bulk learning" approach of Shepanski [47], the "high learning rate" approach of Cater [9], and the reinforcement training concept of Williams [58]. New architectures formed by combining backpropagation with other architectures such as Carpenter and Grossberg's Adaptive Resonance [6,7] and Fukushima's Neocognitron [14,15] can also be anticipated.

Methods for adding weight-dependent terms to a network's error function that create a force pushing the network towards a smaller number of non-zero weights, a smaller number of processing elements, or some other "goal" have also been demonstrated [51,43]. As these on-line methods and related off-line methods are developed further, it may become possible for backpropagation architectures to become "self-adjusting". Statistical tests for determining whether a backpropagation network is "overfitting" or "underfitting" a particular training set may well emerge. These may someday allow the development of automated systems that will optimally craft a backpropagation network for a particular data set.

## 6 Conclusions

As demonstrated by the above results, backpropagation is a new tool for approximating functions on the basis of examples. White [53] has summarized the value of this neural network architecture by noting that "... researchers have been searching a long time for an effective set of approximation functions that adapt their functional form to that of the function to be approximated. By employing hierarchical sigmoided adaptive affine combinations, backpropagation achieves this goal by creating a coalition of functions that have their 'frequencies', 'phases' and 'amplitudes' adjusted to precisely fit the function being approximated." The concept of approximating arbitrary  $L_2$  functions using functional forms that do not depend upon either orthogonality or linear superposition may well turn out to be an important theme not only in neurocomputing, but in mathematics, engineering, physics as well; and possibly even in neuroscience (see the Appendix below for a speculation in this direction). The most important results about backpropagation are undoubtedly yet to come.

After presenting this paper (including Theorem 1) at the 1988 INNS Meeting in September 1988, it was brought to my attention that White and his colleagues had, almost simultaneously, independently discovered a similar theorem [24]. Coincidentally, it was White's earlier paper with Gallant [16] that led me to the idea for the above theorem. Other papers examining the issue of Theorem 1 that have subsequently been brought to my attention include a paper also presented at INNS-88 by Moore and Poggio [37] and the 1987 doctoral dissertation of le Cun [35]. Recent work by Irie and Miyake [26], le Cun [34], and Becker and le Cun [3] is also relevant. The author wishes to thank the referees for numerous helpful comments and suggestions.

This work was partially supported by the SDI Innovative Science and Technology Program under U. S. Army Strategic Defense Command contract DASGGO-88-C-0112 to HNC, Inc.

## Appendix

Until now it has been extremely difficult to believe that the traditional backpropagation neural network architecture was relevant to neurophysiology at the cellular level (this is not true for non-traditional variants of backpropagation such as those of Parker [41,40]). This difficulty follows from the fact that past constructions of the traditional backpropagation architecture have involved non-local processing – which is believed to be impossible in neural tissue. Since the new architecture presented in Section 2 above eliminates the non-locality of backpropagation and makes it a legitimate neural network, while retaining the traditional mathematical form, it may be sensible to reexamine the possible biological relevance of backpropagation. As a start in this direction this Appendix presents a plausible, but highly speculative, hypothetical neurophysiological implementation of backpropagation. It is important to point out that the neurons involved in this proposed neural circuit are almost certainly also involved in other circuits as well, some of which might well be active at the same time. Thus, this hypothesis does not attempt to account for the totality of cortical function, merely one aspect of it – namely, the learning of associations or mappings between nearby cortical regions. The hypothesis is attractive in that, as demonstrated in Section 4 above, multilayer networks can learn virtually any desired associative mapping, unlike most of the simple linear associative schemes that have been proposed. (Although, in fairness, some basically linear schemes are more capable than they might initially appear, in particular: the sparse associative memory of Willshaw, et al. [61] as well as the Bidirectional Associative Memory (BAM) of Kosko [32,33] and the Hopfield network [23], at least when the capacity improvement modifications proposed by Haines and Hecht-Nielsen [19] are incorporated.)

The hypothesis presented here is that backpropagation is used in the cerebral cortex for learning complicated mappings between nearby areas of cortex that are interconnected by the axons of shallow pyramid neurons. As is well known [13,28], the white matter axons of small to medium shallow pyramids are short and only go to nearby areas of cortex as elements of the short association tracts [11]. The axon collaterals of larger, deeper pyramids also go to these adjacent cortical regions. But, unlike the shallow pyramids, other collaterals of the deep pyramids make up the major brain fascicles that interconnect distant cortical regions and send signals to extracortical areas.

The backpropagation hypothesis assumes that the forward pass of the network is active almost all of the time. The backward pass is triggered only occasionally – namely, when it is necessary to learn something. It is assumed that this learning operation is triggered by some sort of "mismatch detection", attention direction "searchlight", or "important event detection" function carried out by thalamic tissue (including the LGN, MGN, pulvinar, and the thalamus proper) and/or the thalamic reticular complex as in the theories of Grossberg [17,18] and Crick [10]. It is assumed that input to cortex from thalamus triggers the backward pass (which then proceeds using Equation [34]). The thalamic input is assumed to modify the behavior of all of the cells of the effected area of cortex, except the deep pyramids – which are exempted, perhaps by means of the action of horizontal cells (which are known to receive thalamic signals and which preferentially synapse with the apical dendrites of the deep pyramids) [11].

The details of the hypothesis are presented in Figure 5. The deep pyramid cells (which are known to carry out an "integrative" function) are assumed to carry out the feed-forward function of the suns – namely, summing the inputs from the shallow pyramids (which, together with normal stellate cells, are assumed to carry out both the forward and backward pass functions of the planets) and applying a sigmoid function to the sum. Thus,

each deep pyramid and the shallow pyramids that feed it make up a sun and planet "unit" of Figures 1 and 2. The feedback function of the suns (summing the back-propagated error signals) is assumed to be carried out by cortical "basket" or "basket stellate" cells which synapse with shallow pyramids over a considerable local area (unlike the normal stellates, which only synapse with a small number of shallow pyramids that are very close) [27]. These basket stellates are assumed to receive input preferentially from the cortico-cortical white matter axons of shallow pyramids of nearby cortical areas (i.e., those participating in the short association tracts of the white matter), but not deep pyramid inputs. The shallow pyramid/normal stellate "planet units" are assumed to receive inputs from deep pyramids exclusively. The supposed operation of this neural circuit is now described. Throughout the discussion it is assumed that the transmitted signals are pulse frequency modulated transmissions ranging continuously between zero firing frequency and some maximum firing frequency (corresponding to the asymptotic 0 and 1 outputs of the backpropagation sigmoid). The need to accommodate negative signals for error feedback is ignored here – presumably some sort of basket stellate or normal stellate/shallow pyramid offset bias could fix this problem.

On the forward pass (i.e., the normal operational association mode) shallow pyramids that have become activated (either by feedforward input from "lower layers", or by external input) send their outputs to a deep pyramid that is nearby. Output axons from these shallow pyramids also leave cortex and travel to nearby cortical regions, where they synapse with basket stellates. However, the basket stellates (unlike the normal stellates) are assumed to be inactive until the thalamic learning signal is present. The output of the deep pyramid that sums and signoids the output of local shallow pyramids is then also sent to nearby cortical areas via the same association bundles containing the shallow pyramid axon collaterals. The deep pyramid axons then synapse with normal stellate/shallow pyramid planet units; thus providing the required forward pass input to the "next layer".

When the thalamus activates the backward pass (which is assumed to only last a short time, although if the thalamic activation is triggered by the hippocampus during sleep this need not be true) the deep pyramids are somehow exempted from interruption of their function. In fact, they keep firing at the same rate they were operating at immediately before the thalamic input occurred. This allows the planet units that receive these inputs to use them in updating their weights by multiplying the incoming deep pyramid signal times the error signal coming in from the basket stellate and adding this to the existing weight. While it is not easy to envision how a couple of cells can carry out these calculations, it does not seem beyond possibility. The shallow pyramids then transmit the product of the basket stellate input (the error for this group of planets) and their weight (either before or after updating). These error signals are then transmitted to the appropriate basket stellate(s) of the next "lower" layer. Note that in accordance with Figures 1 and 2 that these connections must be very specific (unlike the deep pyramid outputs, which must be broadly distributed). In particular, the shallow pyramid white matter axons that participate in implementation of backpropagation *must* functionally connect only to the basket cell(s) associated with the shallow pyramids that activate the deep pyramid that feeds that particular shallow pyramid. This does not necessarily mean that each shallow pyramid has only one or a few collaterals. For example, it is possible that the signals in collaterals that target cells other than those few basket stellates needed by backpropagation may be statistically meaningless because they are randomly uncorrelated with activity in the target region. These connections might be used to implement other networks at other times. In general, we can conclude that for this hypothesis to be correct there must be many more shallow pyramids than deep pyramids and that deep pyramids must have more numerous and more broadly distributed axon collaterals. Clearly, this jibes with the neurophysiological facts. Perhaps this (admittedly crude) cortical backpropagation hypothesis can serve to stimulate some useful thought and discussion.

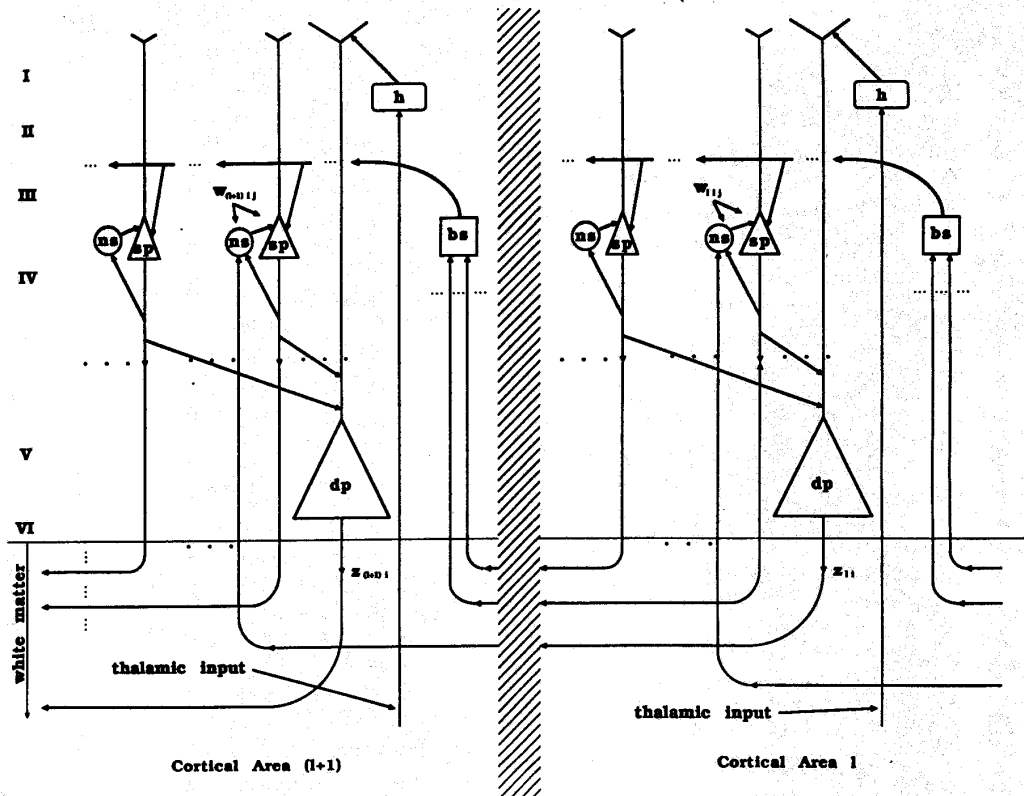


Figure 5: A possible cerebral neurophysiological implementation of backpropagation (essentially Figure 2 done with neurons). In this hypothesis, the deep pyramids (labeled *dp*) carry out the feedforward functions of the suns. Basket stellates (labeled *bs*) carry out the feedback error summation of the suns. Normal stellates (labeled *ns*) and shallow pyramids (labeled *sp*), working together, carry out the functions of the planets (deep pyramid input weighting and weight modification). Thalamic "EEG" signals trigger the backward pass (via collaterals not shown) and also activate horizontal cells (labeled *h*) – which allow the deep pyramids to continue firing during learning by exempting them from the influence of the thalamic signals.

## References

- [1] Anderson, James A., and Rosenfeld, Edward, [Eds.], **Neurocomputing: Foundations of Research**, MIT Press, Cambridge, Massachusetts, 1988.
- [2] Ash, Timur, "Dynamic Node Creation in Backpropagation Networks", Department of Computer Science and Engineering, University of California at San Diego, Preliminary Manuscript, January 1989.
- [3] Becker, Sue, and le Cun Yann, "Improving the Convergence of Back-Propagation Learning with Second Order Methods", Technical Report CRG-TR-88-5, Connectionist Research Group, University of Toronto, Canada, September 1988.
- [4] Bronshtein, I.N., and Semendyayev, K.A., **Handbook of Mathematics**, Third Edition, Van Nostrand Reinhold, New York, 1985.
- [5] Bryson, A.E., and Ho, Yu-Chi, **Applied Optimal Control**, Blaisdell, New York, 1969.
- [6] Carpenter, Gail A., and Grossberg, Stephen, "A Massively Parallel Architecture for a Self-organizing Neural Pattern Recognition Machine", *Computer Vision, Graphics and Image Processing*, **37**, 54-115, 1987.
- [7] Carpenter, Gail A., and Grossberg, Stephen, "ART 2: self-organization of stable category recognition codes for analog input patterns", *Applied Optics*, **26**, No.23, 4919-4930, 1 December 1987.
- [8] Carpenter, Gail A., and Grossberg, Stephen, "Associative learning, adaptive pattern recognition, and cooperative-competitive decision making by neural networks", in Szu, Harold (Ed.), **Optical and Hybrid Computing**, SPIE Institute Series, published as: *SPIE Proc.*, **634**, 218-247, 1986.
- [9] Cater, John, P., "Successfully Using Peak Learning Rates of 10 (and Greater) in Back-Propagation Networks with the Heuristic Learning Algorithm", *Proc. 1987 IEEE International Conference on Neural Networks*, II(645-651), IEEE Press, New York, 1987.
- [10] Crick, Francis H.C., "Function of the thalamic reticular complex: the searchlight hypothesis", *Proc. Nat. Acad. Sci.*, **81**, 4586-4590, 1984. { Note: also published in Anderson, James A., and Rosenfeld, Edward, [Eds.], **Neurocomputing: Foundations of Research**, MIT Press, Cambridge, Massachusetts, 1988. }
- [11] Diamond, Marian C., Scheibel, Arnold B., and Elson, Lawrence M., **The Human Brain Coloring Book**, Barnes & Noble Books, New York, 1985.
- [12] Dunford, Nelson, and Schwartz, Jacob T., **Linear Operators, Part I**, Third Printing, Wiley Interscience Publishers, New York, 1966.
- [13] Feldman, Martin L., **Morphology of the Neocortical Pyramidal Neuron** in Peters, Alan, and Jones, Edward G. [Eds.] **Cerebral Cortex, Volume 1: Cellular Components of the Cerebral Cortex**, 123-200, Plenum Press, New York, 1984.
- [14] Fukushima, Kunihiro, "Neocognitron: A Hierarchical Neural Network Capable of Visual Pattern Recognition", *Neural Networks*, **1**, 119-130, 1988.
- [15] Fukushima, Kunihiro, and Miyake, Sei, "Neocognitron: A New Algorithm for Pattern Recognition Tolerant of Deformations and Shifts in Position", *Pattern Recognition*, **15**, No. 6, 455-469, 1984.
- [16] Gallant, A. Ronald, and White, Halbert, "There exists a neural network that does not make avoidable mistakes", *Proc. 1988 IEEE International Conference on Neural Networks*, pp I(657-664), IEEE Press, New York, 1988.
- [17] Grossberg, Stephen [Ed.], **Neural Networks and Natural Intelligence**, MIT Press, Cambridge, 1988.
- [18] Grossberg, Stephen, **Studies of Mind and Brain**, Reidel, Boston, 1982.
- [19] Haines, Karen, and Hecht-Nielsen, Robert, "A BAM With Increased Information Storage Capacity", *Proc. 1988 IEEE International Conference on Neural Networks*, I(181-190), IEEE Press, New York, 1988.
- [20] Hecht-Nielsen, Robert, "Kolmogorov's Mapping Neural Network Existence Theorem", *Proc. 1987 IEEE International Conference on Neural Networks*, IEEE Press, New York, III(11-13), 1987.
- [21] Hinton, Geoffrey E., Sejnowski, Terrence J., and Ackley, David H., **Boltzmann Machines: Constraint Satisfaction Networks that Learn**, Report CMU-CS-84-157, Carnegie-Mellon University, 1984.
- [22] Hirsch, Morris, "Dynamical Systems Review", a tutorial presented at the 1988 IEEE International Conference on Neural Networks, videotape and notes available from: IEEE Press, New York, 1988.
- [23] Hopfield, J. J., "Neurons With Graded Response Have Collective Computational Properties Like Those of Two-state Neurons", *Proc. Natl. Acad. Sci.*, **81**, 3088-3092, May 1984.
- [24] Hornik, Kurt, Stinchcombe, Maxwell, and White, Halbert, "Multilayer Feedforward Networks Are Universal Approximators", Manuscript, Department of Economics, University of California at San Diego, June 1988.
- [25] Hush, D. R., Salas, J. M., "Improving the Learning Rate of Back-Propagation with the Gradient Reuse Algorithm", *Proc. 1988 IEEE International Conference on Neural Networks*, I(441-446), IEEE Press, New York, 1988.
- [26] Irie, Bunpei, and Miyake, Sei, "Capabilities of Three-layered Perceptrons", *Proc. 1988 IEEE International Conference on Neural Networks*, I(641-648), IEEE Press, New York, 1988.
- [27] Jones, Edward G. and Hendry, Stewart H.C., **Basket Cells** in Peters, Alan, and Jones, Edward G. [Eds.] **Cerebral Cortex, Volume 1: Cellular Components of the Cerebral Cortex**, 309-336, Plenum Press, New York, 1984.
- [28] Jones, Edward G., **Laminar Distribution of Cortical Efferent Cells** in Peters, Alan, and Jones, Edward G. [Eds.] **Cerebral Cortex, Volume 1: Cellular Components of the Cerebral Cortex**, 521-553, Plenum Press, New York, 1984.
- [29] Kohonen, Teuvo, **Self-Organization and Associative Memory**, Second Edition, Springer-Verlag, New York, 1988.
- [30] Kolmogorov, Andrei Nikolaevich, "On the Representation of Continuous Functions of Many Variables by Superposition of Continuous Functions of One Variable and Addition", *Dokl. Akad. Nauk USSR*, **114**, 953-956, 1957.

- [31] Korn, Granino A., and Korn, Theresa M., **Mathematical Handbook for Scientists and Engineers**, Second Edition, McGraw-Hill, New York, 1968.
- [32] Kosko, Bart, "Bidirectional Associative Memories", *IEEE Trans. on Systems, Man, and Cyber.*, **18**, No. 1, 49-60, January/February 1988.
- [33] Kosko, Bart, "Adaptive Bidirectional Associative Memories", *Applied Optics*, **26**, No. 23, 4947-4960, 1 December 1987.
- [34] le Cun, Yann, "A Theoretical Framework for Back-Propagation", Technical Report CRG-TR-88-6, Connectionist Research Group, University of Toronto, Canada, September 1988.
- [35] le Cun, Yann, "Modeles Connexionnistes de l'Apprentissage", Doctoral Dissertation, University of Pierre and Marie Curie, Paris, France, 1987.
- [36] McInerney, John M., Haines, Karen G., Biafore, Steve, and Hecht-Nielsen, Robert, "Can Backpropagation Error Surfaces Have Non-Global Minima?", Department of Electrical and Computer Engineering, University of California at San Diego, Manuscript, August 1988.
- [37] Moore, B. and Poggio, T., "Representation Properties of Multilayer Feedforward Networks", Presented at the 1988 Annual Meeting of the International Neural Network Society, Boston, Massachusetts, 8 Sept 1988. Note: an abstract of this paper was published as a Supplement to Volume 1 of *Neural Networks*, 1988.
- [38] Minsky, Marvin, and Papert, Seymour, **Perceptrons**, Expanded Edition, MIT Press, 1988.
- [39] Parker, David B., "Optimal Algorithms for Adaptive Networks: Second Order Back Propagation, Second Order Direct Propagation, and Second Order Hebbian Learning", *Proc. 1987 IEEE International Conference on Neural Networks*, II(593-600), IEEE Press, New York, 1987.
- [40] Parker, David B., "A Comparison of Algorithms for Neuron-Like Cells", in Denker, John [Ed.], **Proc. Second Annual Conference on Neural Networks for Computing**, Proceedings Vol. 151, 327-332, American Institute of Physics, New York, 1986.
- [41] Parker, David B., "Learning-Logic", Technical Report TR-47, Center for Computational Research in Economics and Management Science, MIT, April 1985.
- [42] Ricotti, L. P., Ragazzini, S., Martinelli, G., "Learning of Word Stress in a Sub-Optimal Second order Back-Propagation Neural Network", *Proc. 1988 IEEE International Conference on Neural Networks*, I-355 - I-361, IEEE Press, New York, 1988.
- [43] Rumelhart, David E., "Parallel Distributed Processing", Plenary Lecture presented at *Proc. 1988 IEEE International Conference on Neural Networks*, San Diego, California, July 1988.
- [44] Rumelhart, David E., Hinton, Geoffrey E., and Williams, Ronald J., "Learning representations by back-propagating errors", *Nature*, **323**, 533-536, 9 October 1986.
- [45] Rumelhart, David E., and McClelland, James L., **Parallel Distributed Processing: Explorations in the Microstructure of Cognition**, Vols. I, II & III, MIT Press, 1986 & 1987.
- [46] Rumelhart, David E., Hinton, Geoffrey E., and Williams, Ronald J., "Learning Internal Representations by Error Propagation", ICS Report 8506, Institute for Cognitive Science, University of California at San Diego, September 1985.
- [47] Shepanski, J. F., "Fast Learning in Artificial Neural Systems: Multilayer Perceptron Training Using Optimal Estimation", *Proc. 1988 IEEE International Conference on Neural Networks*, I(465-472), IEEE Press, New York, 1988.
- [48] Watrous, Raymond, L., "Learning Algorithms for Connectionist Networks: Applied Gradient Methods of Nonlinear Optimization", *Proc. 1987 IEEE International Conference on Neural Networks*, II(619-627), IEEE Press, New York, 1987.
- [49] Werbos, Paul J., "Backpropagation: Past and Future", *Proc. 1988 IEEE International Conference on Neural Networks*, I(343-353), IEEE Press, New York, 1988.
- [50] Werbos, Paul J. "Building and Understanding Adaptive Systems: A Statistical/Numerical Approach to Factory Automation and Brain Research", *IEEE Trans. on Systems, Man, and Cyber.*, SMC-17, No. 1, 7-20, January/February 1987.
- [51] Werbos, Paul J. "Learning How The World Works: Specifications for Predictive Networks in Robots and Brains", *Proc. 1987 IEEE Conf. on Systems, Man, and Cyber.*, IEEE Press, New York, 1987.
- [52] Werbos, Paul J., **Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences**, Ph.D. Thesis, Applied Mathematics, Harvard University, November, 1974.
- [53] White, Halbert, personal communication, 3 October 1988.
- [54] Widrow, Bernard, and Stearns, Samuel D., **Adaptive Signal Processing**, Prentice-Hall, 1985.
- [55] Widrow, Bernard, "Generalization and Information Storage in Networks of ADALINE Neurons", in Yovitts, G. T., **Self-Organizing Systems**, Spartan Books, 1962.
- [56] Widrow, Bernard, and Hoff, Marcian T., Jr., "Adaptive Switching Circuits", *IRE WESCON Conv. Record, Part 4*, 96-104, 1960.
- [57] Williams, Ronald J., "On the Use of Backpropagation in Associative Reinforcement Learning", *Proc. 1988 IEEE International Conference on Neural Networks*, I(623-270), IEEE Press, New York, 1988.
- [58] Williams, Ronald J., "A Class of Gradient-Estimating Algorithms for Reinforcement Learning in Neural Networks", *Proc. 1987 IEEE International Conference on Neural Networks*, II(601-608), IEEE Press, New York, 1987.
- [59] Williams, Ronald J., "Feature Discovery Through Error Correction Learning", ICS Report 8501, Institute of Cognitive Science, University of California at San Diego, May 1985.
- [60] Williams, Ronald J., "Unit Activation Rules For Cognitive Network Models", ICS Report 8303, Institute of Cognitive Science, University of California at San Diego, November 1983.
- [61] Willshaw, D.J., Buneman, O.P. and Longuet-Higgins, H.C., "Non-holographic Associative Memory", *Nature*, **222**, 960-962, June 1969.