

CONNECTIONIST LEARNING PROCEDURES¹

Geoffrey E. Hinton²

June 1987

Technical Report CMU-CS-87-115

**Computer Science Department
Carnegie-Mellon University
Pittsburgh PA 15213**

Abstract

A major goal of research on networks of neuron-like processing units is to discover efficient learning procedures that allow these networks to construct complex internal representations of their environment. The learning procedures must be capable of modifying the connection strengths in such a way that internal units which are not part of the input or output come to represent important features of the task domain. Several interesting gradient-descent procedures have recently been discovered. Each connection computes the derivative, with respect to the connection strength, of a global measure of the error in the the performance of the network. The strength is then adjusted in the direction that decreases the error. These relatively simple, gradient-descent learning procedures work well for small tasks and the new challenge is to find ways of improving the speed of learning so that they can be applied to larger, more realistic tasks.

¹This research was supported by contract N00014-86-K-00167 from the Office of Naval Research and by grant IST-8520359 from the National Science Foundation.

²Author's current address is: Computer Science Department, University of Toronto, 10 Kings College Road, Toronto, Ontario M5S 1A4, Canada.

Table of Contents

1. Introduction	1
2. Connectionist models	1
3. Connectionist research issues	2
3.1. Search	3
3.2. Representation	3
3.3. Learning	4
4. Associative memories without hidden units	5
4.1. Linear associators	5
4.2. Non-linear associative nets	6
4.3. The deficiencies of associators without hidden units	8
5. Simple supervised learning procedures	9
5.1. An LMS procedure for binary threshold units	10
5.2. The perceptron convergence procedure	11
5.3. The deficiencies of simple learning procedures	12
6. Back-propagation: A multilayer LMS procedure	13
6.1. The shape of the error surface	14
6.2. Two examples of back-propagation	14
6.3. Back-propagation as a maximum likelihood procedure	15
6.4. Iterative back-propagation	16
6.5. Self-supervised back-propagation	17
6.6. A reinforcement version of back-propagation	18
6.7. The deficiencies of back-propagation	19
7. Boltzmann machines	20
7.1. Maximizing reinforcement and entropy in a Boltzmann machine	23
8. Maximizing mutual information: A semi-supervised learning procedure	24
9. Unsupervised Hebbian learning	25
9.1. A recent development of unsupervised Hebbian learning	26
10. Competitive learning	26
10.1. The relationship between competitive learning and back-propagation	28
11. Reinforcement learning procedures	30
11.1. Delayed reinforcement	31
11.2. The A_{R,P} procedure	31
11.3. Achieving global optimality by reinforcement learning	32
11.4. The relative payoff procedure	32
11.5. Genetic algorithms	34
11.6. Genetic learning and the relative payoff rule	34
11.7. Iterated genetic hill-climbing	35
12. Discussion	35
12.1. Generalization	35
12.2. The speed of learning	37
12.3. Hardware modularity	38
12.4. Other types of modularity	39
12.5. Conclusion	40

1. Introduction

Recent technological advances in VLSI and computer aided design mean that it is now much easier to build massively parallel machines. This has contributed to a new wave of interest in models of computation that are inspired by neural nets rather than the formal manipulation of symbolic expressions. To understand human abilities like perceptual interpretation, content-addressable memory, commonsense reasoning, and learning it may be necessary to understand how computation is organized in systems like the brain which consist of massive numbers of richly-interconnected but rather slow processing elements.

This paper focusses on the question of how internal representations can be learned in "connectionist" networks. These are a recent subclass of neural net models that emphasize computational power rather than biological fidelity. They grew out of work on early visual processing and associative memories (Hinton and Anderson, 1981; Feldman and Ballard, 1982; Rumelhart, McClelland et. al., 1986). The paper starts by reviewing the main research issues for connectionist models and then describes some of the earlier work on learning procedures for associative memories and simple pattern recognition devices. These learning procedures cannot generate internal representations: They are limited to forming simple associations between representations that are specified externally. Recent research has led to a variety of more powerful connectionist learning procedures that can discover good internal representations and most of the paper is devoted to a survey of these procedures.

2. Connectionist models

Connectionist models typically consist of many simple, neuron-like processing elements called "units" that interact using weighted connections. Each unit has a "state" or "activity level" that is determined by the input received from other units in the network. There are many possible variations within this general framework. One common, simplifying assumption is that the combined effects of the rest of the network on the j^{th} unit are mediated by a single scalar quantity, x_j . This quantity, which is called the "total input" of unit j , is a *linear* function of the activity levels of the units that provide input to j :

$$x_j = \sum_i y_i w_{ji} - \theta_j \quad (1)$$

where y_i is the state of the i^{th} unit, w_{ji} is the weight on the connection from the i^{th} to the j^{th} unit and θ_j is the threshold of the j^{th} unit. The threshold term can be eliminated by giving every unit an extra input connection whose activity level is fixed at 1. The weight on this special connection is the negative of the threshold, and it can be learned in just the same way as the other weights. This method of implementing

thresholds will be assumed throughout the paper. An external input vector can be supplied to the network by clamping the states of some units or by adding an input term, I_j , to the total input of some units. The state of a unit is typically defined to be a non-linear function of its total input. For units with discrete states, this function typically has value 1 if the total input is positive and value 0 (or -1) otherwise. For units with continuous states one typical non-linear transfer function is the logistic function (shown in figure 1).

$$y_j = \frac{1}{1 + e^{-x_j}} \quad (2)$$

Insert figure 1 about here

All the long-term knowledge in a connectionist network is encoded by where the connections are or by their weights, so learning consists of changing the weights or adding or removing connections. The short-term knowledge of the network is normally encoded by the states of the units, but some models also have fast-changing temporary weights or thresholds that can be used to encode temporary contexts or bindings.

There are two main reasons for investigating connectionist networks. First, these networks resemble the brain much more closely than conventional computers. Even though there are many detailed differences between connectionist units and real neurons, a deeper understanding of the computational properties of connectionist networks may reveal principles that apply to a whole class of devices of this kind, including the brain. Second, connectionist networks are massively parallel, so any computations that can be done with these networks can make good use of parallel hardware.

3. Connectionist research issues

There are three main topics of research on connectionist networks: Search, representation, and learning. This paper focuses on learning, but a very brief introduction to search and representation is necessary in order to understand what learning is intended to produce.

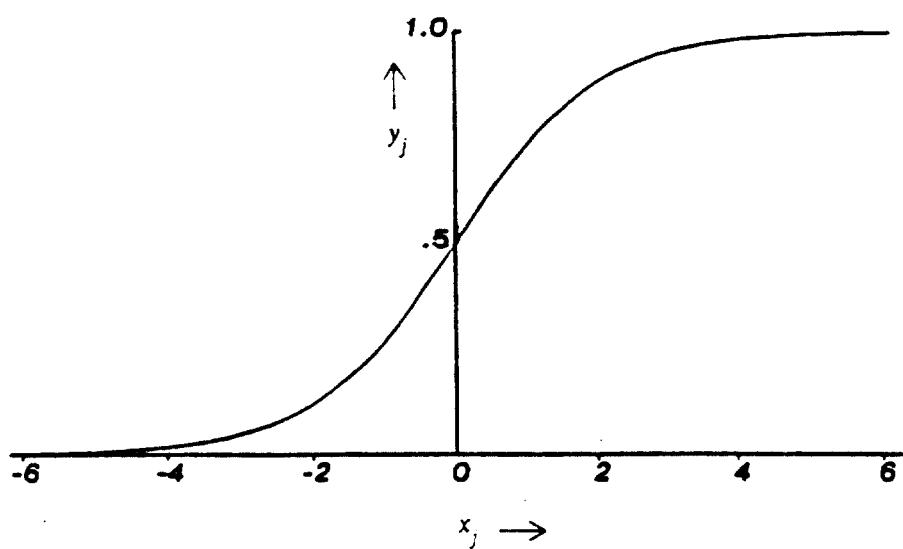


Figure 1: The logistic input-output function defined by equation 2. It is a smoothed version of a step function.

3.1. Search

The task of interpreting the perceptual input, or constructing a plan, or accessing an item in memory from a partial description can be viewed as a constraint-satisfaction search in which information about the current case (i.e. the perceptual input or the partial description) must be combined with knowledge of the domain to produce a solution that fits both these sources of constraint as well as possible (Ballard, Hinton, and Sejnowski, 1983). If each unit represents a piece of a possible solution, the weights on the connections between units can encode the degree of consistency between various pieces. In interpreting an image, for example, a unit might stand for a piece of surface at a particular depth and surface orientation. Knowledge that surfaces usually vary smoothly in depth and orientation can be encoded by positive weights between units that represent nearby pieces of surface at similar depths and similar surface orientations, and negative weights between nearby pieces of surface at very different depths or orientations. The network can perform a search for the most plausible interpretation of the input by iteratively updating the states of the units until they reach a stable state in which the pieces of the solution fit well with each other and with the input. Any one constraint can typically be overridden by combinations of other constraints and this makes the search procedure robust in the presence of noisy data, noisy hardware, or minor inconsistencies in the knowledge.

There are, of course, many complexities: Under what conditions will the network settle to a stable solution? Will this solution be the optimal one? How long will it take to settle? What is the precise relationship between weights and probabilities? These issues are examined in detail by Hummel and Zucker (1983), Hinton and Sejnowski (1983), Geman and Geman (1984), Hopfield and Tank (1985) and Marronquin (1985).

3.2. Representation

For tasks like low-level vision, it is usually fairly simple to decide how to use the units to represent the important features of the task domain. Even so, there are some important choices about whether to represent a physical quantity (like the depth at a point in the image) by the state of a single continuous unit, or by the activities in a set of units each of which indicates its confidence that the depth lies within a certain interval (Ballard, 1986).

The issues become much more complicated when we consider how a complex, articulated structure like a plan or the meaning of a sentence might be represented in a network of simple units.

Some preliminary work has been done on the representation of inheritance hierarchies and the representation of frame-like structures in which a whole object is composed of a number of parts each of which plays a different role within the whole (Minsky, 1977; Hinton, 1981). A recurring issue is the distinction between local and distributed representations. In a local representation, each concept is represented by a single unit (Barlow, 1972; Feldman, 1986). In a distributed representation, the kinds of concepts that we have words for are represented by patterns of activity distributed over many units, and each unit takes part in many such patterns (Hinton, McClelland and Rumelhart, 1986). Distributed representations are usually more efficient than local ones in addition to being more damage resistant. Also, if the distributed representation allows the weights to capture important underlying regularities in the task domain, it can lead to much better generalization than a local representation (Rumelhart and McClelland, 1986; Rumelhart, Hinton and Williams, 1986a). However, distributed representations can make it difficult to represent several different things at the same time and so to use them effectively for representing structures that have many parts playing different roles it may be necessary to have a separate group of units for each role so that the assignment of a filler to a role is represented by a distributed pattern of activity over a group of "role-specific" units.

Much confusion has been caused by the failure to realize that the words "local" and "distributed" refer to the *relationship* between the terms of some descriptive language and a connectionist implementation. If an entity that is described by a single term in the language is represented by a pattern of activity over many units in the connectionist system, and if each of these units is involved in representing other entities, then the representation is distributed. But it is always possible to invent a new descriptive language such that, relative to this language, the very same connectionist system is using local representations.

3.3. Learning

In a network that uses local representations it is often feasible to set all the weights by hand because each weight typically corresponds to a meaningful relationship between entities in the domain. If, however, the network uses distributed representations it may be very hard to program by hand and so a learning procedure may be essential. Some learning procedures, like the perceptron convergence procedure (Rosenblatt, 1962), are only applicable if the desired states of all the units in the network are already specified. This makes the learning task relatively easy. Other, more recent, learning procedures operate in networks that contain "hidden" units (Hinton, and Sejnowski, 1986) whose desired states are

not specified (either directly or indirectly) by the input or the desired output of the network. This makes learning much harder because the learning procedure must (implicitly) decide what the hidden units should represent. The learning procedure is therefore constructing new representations and the results of learning can be viewed as a numerical solution to the problem of whether to use local or distributed representations.

Connectionist learning procedures can be divided into three broad classes: Supervised procedures which require a teacher to specify the desired output vector, reinforcement procedures which only require a single scalar evaluation of the output, and unsupervised procedures which construct internal models that capture regularities in their input vectors without receiving any additional information. As we shall see, there are often ways of converting one kind of learning procedure into another.

4. Associative memories without hidden units

Several simple kinds of connectionist learning have been used extensively for storing knowledge in simple associative networks which consist of a set of input units that are directly connected to a set of output units. Since these networks do not contain any hidden units, the difficult problem of deciding what the hidden units should represent does not arise. The aim is to store a set of associations between input vectors and output vectors by modifying the weights on the connections. The representation of each association is typically distributed over many connections and each connection is involved in storing many associations. This makes the network robust against minor physical damage and it also means that weights tend to capture regularities in the set of input-output pairings, so the network tends to generalize these regularities to new input vectors that it has not been trained on (Anderson and Hinton, 1981).

4.1. Linear associators

In a linear associator, the state of an output unit is a linear function of the total input that it receives from the input units (see Eq. 1). A simple, Hebbian procedure for storing a new association (or "case") is to increment each weight, w_{ji} , between the i^{th} input unit and the j^{th} output unit by

$$\Delta w_{ji} = y_i y_j \quad (3)$$

where y_i and y_j are the activities of an input and an output unit. After a set of associations have been stored, the weights encode the cross-correlation matrix between the input and output vectors. If the input vectors are orthogonal and have length 1, the associative memory will exhibit perfect recall. Even though each weight is involved in storing many different associations, each input vector will produce

exactly the correct output vector (Kohonen, 1977).

If the input vectors are not orthogonal, the simple Hebbian storage procedure is not optimal. For a given network and a given set of associations, it may be impossible to store all the associations perfectly, but we would still like the storage procedure to produce a set of weights that minimizes some sensible measure of the differences between the desired output vectors and the vectors actually produced by the network. This "error measure" can be defined as

$$E = \frac{1}{2} \sum_{j,c} (y_{j,c} - d_{j,c})^2$$

where $y_{j,c}$ is the actual state of output unit j in input-output case c , and $d_{j,c}$ is its desired state. Kohonen (1977) shows that the weight matrix, W , that minimizes this error measure can be expressed analytically as:

$$W = BA^*$$

where B is the matrix whose columns are the desired output vectors and A^* is the pseudo-inverse of the matrix A whose columns are the input vectors. A simple way to compute W is to use an iterative storage procedure that repeatedly sweeps through the whole set of associations and modifies each weight by a small amount in the direction that reduces the error measure. This is a version of the LMS learning procedure described in section 5. The cost of finding an optimal set of weights (in the least squares sense of optimal) is that storage ceases to be a simple "one-shot" process. To store one new association it is necessary to sweep through the whole set of associations many times.

4.2. Non-linear associative nets

Willshaw (1981) described an associative net in which both the units and the weights have just two states: 1 and 0. The weights all start at 0, and associations are stored by setting a weight to 1 if ever its input and output units are both on in any association (see figure 2). To recall an association, each output unit must have its threshold dynamically set to be just less than m , the number of active input units. If the output unit should be on, the m weights coming from the active input units will have been set to 1 during storage, so the output unit is guaranteed to come on. If the output unit should be off, the probability of erroneously coming on is given by the probability that all m of the relevant weights will have been set to 1 when storing other associations. Willshaw showed that associative nets can make efficient use of the information capacity of the weights. If the number of active input units is the log of the total number of input units, the probability of incorrectly activating an output unit can be made very low even when the

network is storing close to 0.69 of its information theoretic capacity.

.....
Insert figure 2 about here
.....

An associative net in which the input units are identical with the output units can be used to associate vectors with themselves. This allows the network to complete a partially specified input vector. If the input vector is a very degraded version of one of the stored vectors, it may be necessary to use an iterative retrieval process. The initial states of the units represent the partially specified vector, and the states of the units are then updated many times until they settle on one of the stored vectors. Theoretically, the network could oscillate, but Willshaw (1971) and Anderson and Mozer (1981) showed that iterative retrieval normally works well. Hopfield (1982) showed that if the units are updated one at a time the iterative retrieval process can be viewed as a form of gradient descent in an "energy function". Hopfield nets store vectors whose components are all +1 or -1 using the simple storage procedure described in equation 3. To retrieve a stored vector from a partial description (which is a vector containing some 0 components), we start the network at the state specified by the partial description and then repeatedly update the states of units one at a time. The units can be chosen in random order or in any other order provided each unit is chosen finitely often. Hopfield observed that the behavior of the network is governed by the global energy function³

$$E = -\frac{1}{2} \sum_{i,j} s_i s_j w_{ij} + \sum_j s_j \theta_j \quad (4)$$

where s_i and s_j are the states of two units. Each time a unit updates its state, it adopts the state that minimizes this energy function because the decision rule used to update a unit is simply the derivative of the energy function. The unit adopts the state +1 if its "energy-gap" is positive and the state -1 otherwise, where the energy gap is defined as

$$\Delta E_j = E(s_j = -1) - E(s_j = +1) = \sum_i s_i w_{ij} - \theta_j \quad (5)$$

So the energy must decrease until the network settles into a local minimum of the energy function. We can therefore view the retrieval process in the following way: The weights define an "energy landscape" over global states of the network and the stored vectors are local minima in this landscape. The retrieval process consists of moving downhill from a starting point to a nearby local minimum.

³The energy function should not be confused with the error function described earlier. Gradient descent in the energy function is performed by changing the states of the units, not the weights.

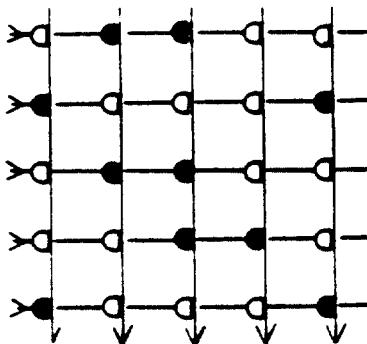


Figure 2: An associative net (Willshaw, 1981). The input vector comes in at the left and the output vector comes out at the bottom (after thresholding). The solid weights have value 1 and the open weights have value 0. The network is shown after it has stored the associations $01001 \rightarrow 10001$, $10100 \rightarrow 01100$, $00010 \rightarrow 00010$.

If too many vectors are stored, there may be spurious local minima caused by interactions between the stored vectors. Also, the basins of attraction around the correct minima may be long and narrow instead of round, so a downhill path from a random starting point may not lead to the nearest local minimum. These problems can be alleviated by using a process called "unlearning" (Hopfield, Feinstein and Palmer, 1983; Crick and Mitchison, 1983).

For random n-component vectors, the storage capacity of a Hopfield net is only $1/k\log(n)$ bits per weight, even though the weights are integers with $c+1$ different values, where c is the number of cases stored. This is much worse than Willshaw's associative net. It can be improved by introducing thresholds (Weisbuch and Fogelman-Soulie, 1985). It can be increased much more by sacrificing the simple one-shot storage procedure and explicitly training the network on typical noisy retrieval tasks using the threshold LMS or perceptron convergence procedures described below.

4.3. The deficiencies of associators without hidden units

If the input vectors are orthogonal, or if they are made to be close to orthogonal by using high-dimensional random vectors (as in a Hopfield net), associators with no hidden units perform well using a simple Hebbian storage procedure. If the set of input vectors satisfy the much weaker condition of being linearly independent, associators with no hidden units can learn to give the correct outputs provided an iterative learning procedure is used. Unfortunately, linear independence does not hold for most tasks that can be characterized as mapping input vectors to output vectors because the number of relevant input vectors is typically much larger than the number of components in each input vector. The required mapping typically has a complicated structure that can only be expressed using multiple layers of hidden units.⁴ Consider, for example, the task of identifying an object when the input vector is an intensity array and the output vector has a separate component for each possible name. If a given type of object can be either black or white, the intensity of an individual pixel (which is what an input unit encodes) cannot provide any direct evidence for the presence or absence of an object of that type. So the object cannot be identified by using weights on direct connections from input to output units. Obviously, it is necessary to explicitly extract relationships among intensity values (such as edges) before trying to identify the object. Actually, extracting edges is just a small part of the problem. If recognition is to have the generative capacity to handle novel images of familiar objects the network must somehow encode the

⁴It is always possible to redefine the units and the connectivity so that multiple layers of simple units become a single layer of much more complicated units. But this redefinition does not make the problem go away.

systematic effects of variations in lighting and viewpoint, partial occlusion by other objects, and deformations of the object itself. There is a tremendous gap between these complex regularities and the regularities that can be captured by an associative net that lacks hidden units.

5. Simple supervised learning procedures

Consider a network that has input units which are directly connected to output units whose states (i.e. activity levels) are a continuous smooth function of their total input. Suppose that we want to train the network to produce particular "desired" states of the output units for each member of a set of input vectors. A measure of how poorly the network is performing with its current set of weights is:

$$E = \frac{1}{2} \sum_{j,c} (y_{j,c} - d_{j,c})^2 \quad (6)$$

where $y_{j,c}$ is the actual state of output unit j in input-output case c , and $d_{j,c}$ is its desired state.

We can minimize the error measure given in Eq. 6 by starting with any set of weights and repeatedly changing each weight by an amount proportional to $\partial E / \partial w$.

$$\Delta w_{ji} = -\epsilon \frac{\partial E}{\partial w_{ji}} \quad (7)$$

In the limit, as ϵ tends to 0 and the number of updates tends to infinity, this learning procedure is guaranteed to find the set of weights that gives the least mean squared (LMS) error. The value of $\partial E / \partial w$ is obtained by differentiating Eq. 6 and Eq. 1.

$$\frac{\partial E}{\partial w_{ji}} = \sum_{\text{cases}} \frac{\partial E}{\partial y_j} \cdot \frac{dy_j}{dx_j} \cdot \frac{\partial x_j}{\partial w_{ji}} = \sum_{\text{cases}} (y_j - d_j) \cdot \frac{dy_j}{dx_j} \cdot y_i \quad (8)$$

If the output units are linear, the term dy_j/dx_j is a constant.

The LMS learning procedure has a simple geometric interpretation if we construct a multi-dimensional "weight-space" that has an axis for each weight and one extra axis (called "height") that corresponds to the error measure. For each combination of weights, the network will have a certain error which can be represented by the height of a point in weight space. These points form a surface called the "error-surface". For networks with linear output units and no hidden units, the error surface always forms a bowl whose horizontal cross-sections are ellipses and whose vertical cross-sections are parabolas. Since the bowl only has one minimum,⁵ gradient descent on the error-surface is guaranteed to find it. If the output units have a non-linear but monotonic transfer function, the bowl is deformed but

⁵This minimum may be a whole subspace.

still only has one minimum, so gradient descent still works. However, it will be very slow at points in weight-space where the gradient of the transfer function approaches zero for the output units that are in error.

The "batch" version of the LMS procedure sweeps though all the cases accumulating $\partial E / \partial w$ before changing the weights, and so it is guaranteed to move in the direction of steepest descent. The "on-line" version, which requires less memory, updates the weights after each input-output case (Widrow and Hoff, 1960)⁶. This may sometimes increase the total error, E, but by making the weight changes sufficiently small the total change in the weights after a complete sweep through all the cases can be made to approximate steepest descent arbitrarily closely.

5.1. An LMS procedure for binary threshold units

Binary threshold units use a step function, so the term dy_j/dx_i is infinite at the threshold and zero elsewhere and the LMS procedure must be modified to be applicable to these units. One method is to redefine the error of an output unit, j, for a given input case to be

$$E_{j,c}^* = \begin{cases} 0 & \text{if the output unit has the right state,} \\ \frac{1}{2}x_{j,c}^2 & \text{if the output unit has the wrong state.} \end{cases}$$

The derivative of the new error measure with respect to $x_{j,c}$ then becomes

$$\frac{\partial E_{j,c}^*}{\partial x_{j,c}} = \begin{cases} 0 & \text{when the output unit behaves correctly,} \\ -x_{j,c} & \text{when the output unit behaves incorrectly} \end{cases}$$

So the "threshold LMS procedure" becomes:

$$\Delta w_{ji} = -\epsilon \sum_c \frac{\partial E_{j,c}^*}{\partial x_{j,c}} \cdot y_{i,c}$$

⁶The term "LMS" is usually reserved for the on-line version

5.2. The perceptron convergence procedure

The usual form of the perceptron convergence procedure is related to the on-line version of the threshold LMS procedure in the following way: The magnitude of $\partial E^*_{j,c} / \partial x_{j,c}$ is ignored and only its sign is taken into consideration. So the weight changes are:

$$\Delta w_{j,c} = \begin{cases} 0 & \text{when the output unit behaves correctly,} \\ +\epsilon y_{i,c} & \text{if the unit is off and should be on,} \\ -\epsilon y_{i,c} & \text{if the unit is on and should be off.} \end{cases}$$

Because it ignores the magnitude of $\partial E^*_{j,c} / \partial x_{j,c}$ this procedure does not even stochastically approximate steepest descent in the type of error function used in LMS procedures. Even with very small ϵ , it is quite possible for the error to rise after a complete sweep through all the cases. However, each time the weights are updated, the perceptron convergence procedure is guaranteed to reduce the value of a different cost measure that is defined solely in terms of weights.

To picture the LMS procedure we introduced a space with one dimension for each weight and one extra dimension for the mean squared error in the output vectors. To picture the perceptron convergence procedure, we do not need the extra dimension for the error. For simplicity we shall consider a network with only one output unit. Each case corresponds to a constraint hyperplane in weight space. If the weights are on one side of this hyperplane, the output unit will behave correctly and if they are on the other side it will behave incorrectly (see figure 3). To behave correctly for all cases, the weights must lie on the correct side of all the hyperplanes, so the combinations of weights that give perfect performance form a convex set. Any set of weights in this set will be called "ideal".

Insert figure 3 about here

The perceptron convergence procedure considers the constraint planes one at a time, and whenever the current combination of weights is on the wrong side, it moves it perpendicularly towards the plane. This reduces the distance between the current combination of weights and any of the ideal combinations. So provided the weights move by less than twice the distance to the violated constraint plane, a weight update is guaranteed to reduce the measure

$$\sum_i (w_{i,actual} - w_{i,ideal})^2$$

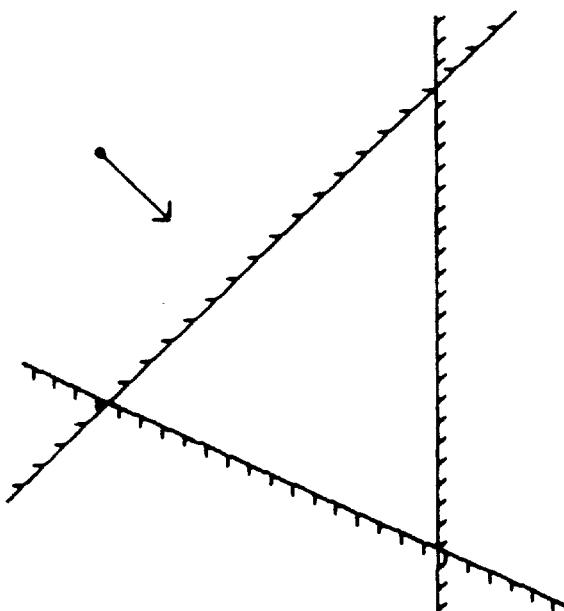


Figure 3: Some hyperplanes in weight-space. Each plane represents the constraint on the weights caused by a particular input-output case. If the weights lie on the correct (unshaded) side of the plane, the output unit will have the correct state for that case. Provided the weight changes are proportional to the activities of the input lines, the perceptron convergence procedure moves the weights perpendicularly towards a violated constraint plane.

The perceptron convergence procedure has many nice properties, but it also has some serious problems. Unlike the threshold LMS procedure, the measure that it is guaranteed to reduce is undefined when there is no set of weights that will do the job perfectly. Also, there are obvious problems in trying to generalize to more complex, multilayered nets in which the ideal combinations of weights do not form a single convex set, because the idea of moving towards the ideal region of weight space breaks down. It is therefore not surprising that the more sophisticated procedures required for multilayer nets are generalizations of the LMS procedure rather than the perceptron convergence procedure: They learn by decreasing a performance error, not a distance in weight-space.

5.3. The deficiencies of simple learning procedures

The major deficiency of both the LMS and perceptron convergence procedures is that most "interesting" mappings between input and output vectors cannot be captured by any combination of weights in such simple networks, so the guarantee that the learning procedure will find the best possible combination of weights is of little value. Consider, for example, a network composed of two input units and one output unit. There is no way of setting the two weights and one threshold to solve the very simple task of producing an output of 1 when the input vector is $(1, 1)$ or $(0, 0)$ and an output of 0 when the input vector is $(1, 0)$ or $(0, 1)$. Minsky and Papert (1969) give a clear analysis of the limitations on what mappings can be computed by three-layered nets. They focus on the question of what pre-processing must be done by the units in the intermediate layer to allow a task to be solved. They generally assume that the pre-processing is fixed, and so they avoid the problem of how to make the units in the intermediate layer learn useful predicates. So, from the learning perspective, their intermediate units are not true hidden units.

Another deficiency of the LMS and perceptron learning procedures is that gradient descent may be very slow if the elliptical cross-section of the error surface is very elongated so that the surface forms a long ravine with steep sides and a very low gradient along the ravine. In this case, the gradient at most points in the space is almost perpendicular to the direction towards the minimum. If the coefficient ϵ in Eq. 7 is large, there are divergent oscillations across the ravine, and if it is small the progress along the ravine is very slow. A standard method for speeding the convergence in such cases is recursive least squares (Widrow and Stearns, 1985). Various other methods have also been suggested (Amari, 1967; Parker, 1987; Plaut and Hinton, 1987).

We now consider learning in more complex networks that contain hidden units. The next five sections describe a variety of supervised, unsupervised, and reinforcement learning procedures for these nets.

6. Back-propagation: A multilayer LMS procedure

The "back-propagation" learning procedure (Rumelhart, Hinton and Williams, 1986a, 1986b) is a generalization of the LMS procedure that works for networks which have layers of hidden units between the input and output units. These multilayer networks can compute much more complicated functions than networks that lack hidden units, but the learning is generally much slower because it must explore the space of possible ways of using the hidden units. There are now many examples in which back-propagation constructs interesting internal representations in the hidden units, and these representations allow the network to generalize in sensible ways. Variants of the procedure were discovered independently by Le Cun (1985) and Parker (1985).

In a multilayer network it is possible, using Eq. 8, to compute $\partial E / \partial w_{ji}$ for all the weights in the network provided we can compute $\partial E / \partial y_j$ for all the units that have modifiable incoming weights. In a system that has no hidden units, this is easy because the only relevant units are the output units, and for them $\partial E / \partial y_j$ is found by differentiating the error function in Eq. 6. But for hidden units, $\partial E / \partial y_j$ is harder to compute. The central idea of back-propagation is that these derivatives can be computed efficiently by starting with the output layer and working backwards through the layers. For each input-output case, c, we first use a forward pass, starting at the input units, to compute the activity levels of all the units in the network. Then we use a backward pass, starting at the output units, to compute $\partial E / \partial y_j$ for all the hidden units. For a hidden unit, j, in layer J the only way it can affect the error is via its effects on the units, k, in the next layer, K (assuming units in one layer only send their outputs to units in the layer above). So we have

$$\frac{\partial E}{\partial y_j} = \sum_k \frac{\partial E}{\partial y_k} \cdot \frac{dy_k}{dx_k} \cdot \frac{dx_k}{dy_j} = \sum_k \frac{\partial E}{\partial y_k} \cdot \frac{dy_k}{dx_k} \cdot w_{kj} \quad (9)$$

where the index c has been suppressed for clarity. So if $\partial E / \partial y_k$ is already known for all units in layer K, it is easy to compute the same quantity for units in layer J. Notice that computation performed during the backward pass is very similar in form to the computation performed during the forward pass (though it propagates error derivatives instead of activity levels, and it is entirely linear in the error derivatives).

6.1. The shape of the error surface

In networks without hidden units, the error surface only has one minimum (provided the units use smooth monotonic transfer functions). With hidden units, the error surface may contain many local minima, so it is possible that steepest descent in weight space will get stuck at poor local minima. In practice, this does not seem to be a serious problem. Back-propagation has been tried for a wide variety of tasks and poor local minima are rarely encountered, provided the network contains a few more units and connections than are required for the task. One reason for this is that there are typically a very large number of qualitatively different perfect solutions, so we avoid the typical combinatorial optimization task in which one minimum is slightly better than a large number of other, widely separated minima.

In practice, the most serious problem is the speed of convergence, not the presence of non-global minima. This is discussed further in section 12.3.

6.2. Two examples of back-propagation

Back-propagation is an effective learning technique when the mapping from input vectors to output vectors contains both regularities and exceptions. For example, in mapping from a string of English letters to a string of English phonemes there are many regularities but there are also exceptions such as the word "women". Sejnowski and Rosenberg (1987) have shown that a network with one hidden layer can be trained to pronounce letters surprisingly well. The input layer encodes the identity of the letter to be pronounced using a different unit for each possible letter, and also the identities of the three previous letters and three following letters in the text (spaces and punctuation are treated as special kinds of letters). The output layer encodes a phoneme using 21 articulatory features and 5 features for stress and syllable boundaries. There are 80 hidden units each of which is connected to all the input and output units. After extensive training, the network generalizes well to new examples which demonstrates that it captures the regularities of the mapping.

A network with three hidden layers was used by Hinton (1986) to learn the family relationships between 24 different people. The input vector specified a person and a relationship and the network was trained to produce the related person as output. The network was trained on 100 of the 104 possible relationships and then tested on the remaining 4. It generalized correctly because during the training it learned to represent each of the people in terms of important features such as age, nationality, and the branch of the family tree that they belonged to. Using these underlying features, much of the information

about family relationships can be captured by a fairly small number of "micro-inferences" between features. For example, the father of a middle-aged person is an old person, and the father of an Italian person is an Italian person. Figures 4 to 6 show the family trees, the architecture of the network and the "receptive fields" of some of the internal units.

Insert figures 4, 5, 6 about here

The input and output encodings used a different unit to represent each person and relationship, so all pairs of people are equally similar in the input and output encoding: The encodings do not give any clues about what the important features are. The learning procedure could only discover these features by searching for a set of features that made it easy to express the associations. Once these features were discovered, the *internal* representation of each person was a distributed pattern of activity and similar people were represented by similar patterns. Thus the network constructs its own internal similarity metric. This is a significant advance over simulations in which good generalization is achieved because the experimenter chooses representations that already have an appropriate similarity metric.

6.3. Back-propagation as a maximum likelihood procedure

If we interpret each output vector as a specification of a conditional probability distribution over a set of output vectors given an input vector, we can interpret the back-propagation learning procedure as a method of finding weights that maximize the likelihood of generating the desired conditional probability distributions. Two examples of this kind of interpretation will be described.

Suppose we only attach meaning to binary output vectors and we treat a real-valued output vector as a way of specifying a probability distribution over binary vectors.⁷ We imagine that a real-valued output vector is stochastically converted into a binary vector by treating the real-values as the probabilities that individual components have value 1, and assuming independence between components. For simplicity, we can assume that the desired vectors used during training are binary vectors, but this is not necessary. Given a set of training cases, it can be shown that the likelihood of producing *exactly* the desired vectors is maximized when we minimize the cross-entropy, C , between the desired and actual conditional probability distributions:

$$C = \sum_{j,c} d_{j,c} \log_2(y_{j,c}) + (1 - d_{j,c}) \log_2(1 - y_{j,c})$$

⁷Both the examples of back-propagation described above fit this interpretation.

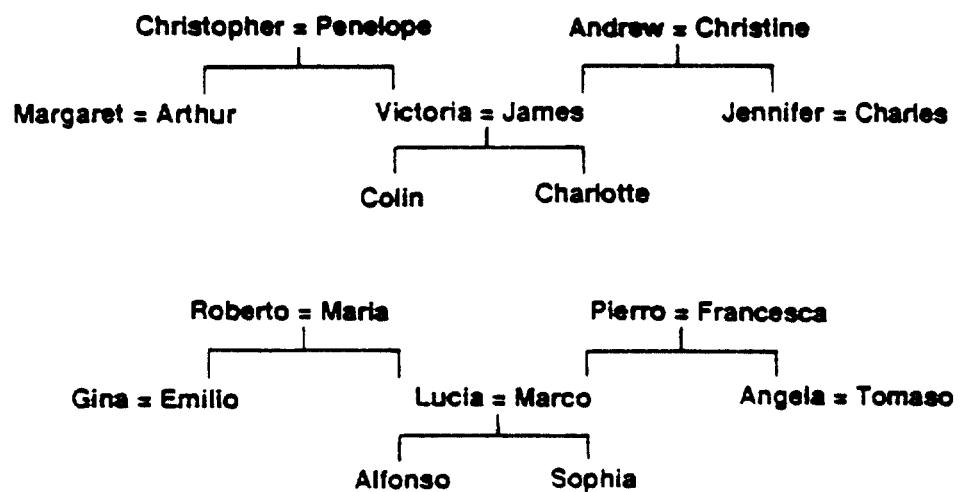


Figure 4: Two isomorphic family trees.

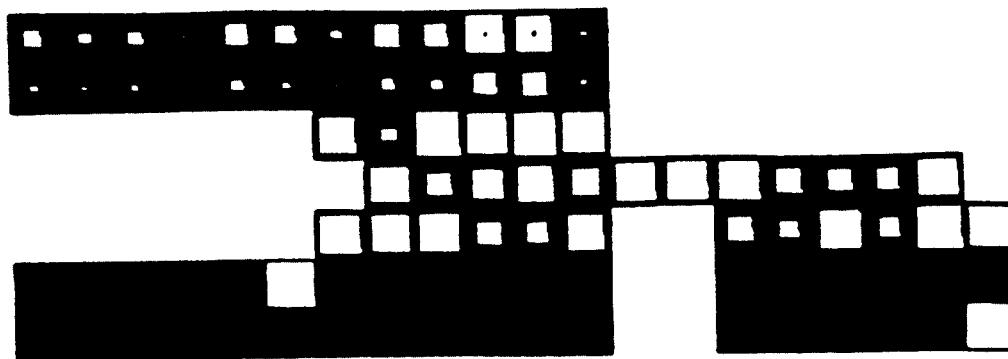


Figure 5: The activity levels in a five layer network after it has learned. The bottom layer has 24 input units on the left for representing person 1 and 12 units on the right for representing the relationship. The white squares inside these two groups show the activity levels of the units. There is one active unit in the first group (representing Colin) and one in the second group (representing has-aunt). Each of the two groups of input units is totally connected to its own group of 6 units in the second layer. These two groups of 6 must learn to encode the input terms as distributed patterns of activity. The second layer is totally connected to the central layer of 12 units, and this layer is connected to the penultimate layer of 6 units. The activity in the penultimate layer must activate the correct output units, each of which stands for a particular person 2. In this case, there are two correct answers (marked by black dots) because Colin has two aunts. Both the input and output units are laid out spatially with the English people in one row and the isomorphic Italians immediately below.

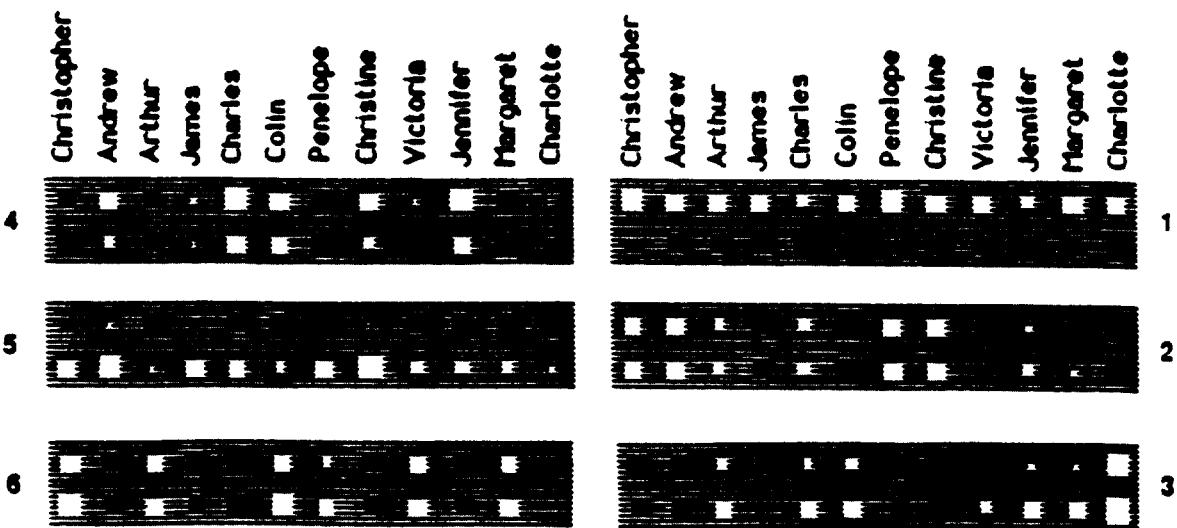


Figure 6: The weights from the 24 input units that represent people to the 6 units in the second layer that learn distributed representations of people. White rectangles stand for excitatory weights, black for inhibitory weights, and the area of the rectangle encodes the magnitude of the weight. The weights from the 12 English people are in the top row of each unit. Beneath each of these weights is the weight from the isomorphic Italian. During the learning, each weight was given a tendency to decay towards zero. This tendency is balanced by the error-gradient, so the final magnitude of a weight indicates how useful it is in reducing the error.

where $d_{j,c}$ is the desired probability of output unit j in case c and $y_{j,c}$ is its actual probability.

So, under this interpretation of the output vectors, we should use the cross-entropy function rather than the squared difference as our cost measure. In practice, this helps to avoid a problem caused by output units which are firmly off when they should be on (or vice versa). These units have a very small value of $\partial y / \partial x$ so they need a large value of $\partial E / \partial y$ in order to change their incoming weights by a reasonable amount. When an output unit that should have an activity level of 1 changes from a level of .0001 to a level of .001, the squared difference from 1 only changes slightly, but the cross-entropy decreases a lot.

This way of interpreting back-propagation raises the issue of whether, under some other interpretation of the output vectors, the squared error might not be the correct measure for performing maximum likelihood estimation. In fact, Richard Golden (1987) has shown that minimizing the squared error is equivalent to maximum likelihood estimation if both the actual and the desired output vectors are treated as the centers of gaussian probability density functions over the space of all real vectors. So the "correct" choice of cost function depends on the way the output vectors are most naturally interpreted.

6.4. Iterative back-propagation

Rumelhart, Hinton, and Williams (1986a) show how the back-propagation procedure can be applied to iterative networks in which there are no limitations on the connectivity. A network in which the states of the units at time t determine the states of the units at time $t+1$ is equivalent to a layered net which has one layer for each time slice. Each weight in the iterative network is implemented by a whole set of identical weights in the corresponding layered net, one for each time slice. In the iterative net, the error is typically the difference between the actual and desired final states of the network, and to compute the error derivatives it is necessary to back-propagate through time, so the history of states of each unit must be stored. Each weight will have many different error derivatives, one for each time step, and the sum of all these derivatives is used to determine the weight change.

.....
Insert figure 7 about here
.....

Back-propagation in iterative nets can be used to train a network to generate sequences or to recognize sequences or to complete sequences. Examples are given in Rumelhart, Hinton and Williams (1986b). Alternatively, it can be used to store a set of patterns by constructing a point attractor for each

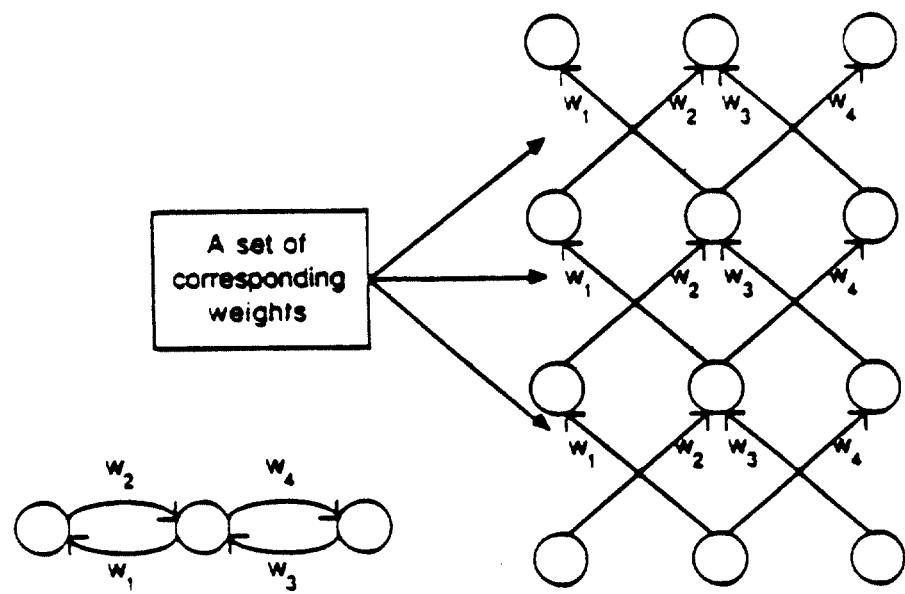


Figure 7: On the left is a simple iterative network that is run synchronously for three iterations. On the right is the equivalent layered network.

pattern. Unlike the simple storage procedure used in a Hopfield net, or the more sophisticated storage procedure used in a Boltzmann machine (see section 7), back-propagation takes into account the path used to reach a point attractor. So it will not construct attractors that cannot be reached from the normal range of starting points on which it is trained.⁸

There is some hope that an iterative net can be trained to recognize phonemes. The speech, in a form such as spectral coefficients, would be presented to the network one time-frame at a time, and the phoneme identities would be output after a short delay. Back-propagation of the errors would cause the hidden units to extract features that were helpful for recognition, and it would also cause some hidden units to "remember" relevant aspects of the earlier input. These hidden units would then form an internal context which would aid recognition of the next phoneme.⁹ In principle, different hidden units could represent many different aspects of the earlier input. This would overcome a major limitation of the hidden markov model approach to speech recognition (Bahl, Jelinek, and Mercer, 1983). Hidden markov models are stochastic finite state automata, and the only internal context available to a finite state automaton is the identity of the current state. To represent 20 bits of information about the past a hidden markov model would need 2^{20} states, and the corresponding recognizer would need a state with up to 2^{20} components. An iterative back-propagation net which had learned to deploy its hidden units to good effect might be able to manage with very much less than 2^{20} units by exploiting patterns in the useful states of the recognition process.

Kevin Lang (1987) has shown that an iterative back-propagation net compares favorably with hidden markov models at the task of distinguishing the words "bee", "dee", "ee", and "vee" spoken by many different male speakers in a noisy environment.

6.5. Self-supervised back-propagation

One drawback of the standard form of back-propagation is that it requires an external supervisor to specify the desired states of the output units. It can be converted into an unsupervised procedure by using the input itself to do the supervision, using a multilayer "encoder" network (Ackley et. al., 1985) in

⁸A back-propagation net that uses asymmetric connections (and synchronous updating) is not guaranteed to settle to a single stable state. To encourage it to construct a point attractor, rather than a limit cycle, the point attractor can be made the desired state for the last few iterations.

⁹It is also possible to use connections that have various different time delays so that units in higher layers can directly "observe" the time-varying behavior of units in lower layers.

which the desired output vector is identical with the input vector. The network must learn to compute an approximation to the identity mapping for all the input vectors in its training set, and if the middle layer of the network contains fewer units than the input layer, the learning procedure must construct a compact, invertible code for each input vector. This code can then be used as the input to later stages of processing.

The use of self-supervised back-propagation to construct compact codes resembles the use of principal components analysis to perform dimensionality reduction, but it has the advantage that it allows the code to be a non-linear transform of the input vector. This form of back-propagation has been used successfully to compress images (Cottrell, personal communication, 1986) and to compress speech waves (Elman & Zipser, 1987). A variation of it has been used to extract the underlying degrees of freedom of simple shapes (Saund, 1986).

It is also possible to use back-propagation to predict one part of the perceptual input from other parts. For example, to predict one patch of an image from neighboring patches it is probably helpful to use the hidden units to explicitly extract edges. In domains with sequential structure, one portion of a sequence can be used as input and the next term in the sequence can be the desired output. This forces the network to extract features that are good predictors. If this is applied to the speech wave, the states of the hidden units will form a non-linear predictive code. It is not yet known whether such codes are more helpful for speech recognition than linear predictive coefficients.

A different variation of self-supervised back-propagation is to insist that all or part of the code in the middle layer change as slowly as possible with time. This can be done by making the desired state of each of the middle units be the state it actually adopted for the previous input vector. This forces the network to use similar codes for input vectors that occur at neighboring times, which is a sensible principle if the input vectors are generated by a process whose underlying parameters change more slowly than the input vectors themselves.

6.6. A reinforcement version of back-propagation

David Rumelhart has shown that back-propagation can be extended to cases in which the desired output vector is unknown, but some transformed "image" of it is known. For example, in trying to imitate a sound, a network might produce an output vector which specifies how to move the speech articulators. Suppose that the network does not receive any direct information about what it should do with its

articulators but it does "know" the desired sound and the actual sound, which is the transformed "image" of the output vector. Rumelhart uses an additional network (which he calls a mental model) that first learns to map from output vectors to their images. Once this mapping has been learned, back-propagation through the mental model can convert error derivatives for the "images" into error derivatives for the output vectors of the basic network.

Munro (1987) has shown that this idea can be applied even when the image of an output vector is simply a single scalar value – the reinforcement. First, the mental model learns to predict expected reinforcement from the combination of the input vector and the output vector. Then the derivative of the expected reinforcement can be back-propagated through the mental model to get the reinforcement derivatives for each component of the output vector of the basic network.

6.7. The deficiencies of back-propagation

Despite its impressive performance on relatively small problems, and its promise as a widely applicable mechanism for extracting the underlying structure of a domain, back-propagation is inadequate, in its current form, for larger tasks because the learning time scales poorly. Empirically, the learning time on a serial machine is very approximately $\text{order}(N^3)$ where N is the number of weights in the network. The time for one forward and one backward pass is $\text{order}(N)$. The number of training examples is typically $\text{order}(N)$, assuming the amount of information per output vector is held constant and enough training cases are used to strain the storage capacity of the network (which is about 2 bits per weight). The number of times the weights must be updated is also approximately $\text{order}(N)$. This is an empirical observation and depends on the nature of the task.¹⁰ On a parallel machine that used a separate processor for each connection, the time would be reduced to approximately $\text{order}(N^2)$. Back-propagation can probably be improved by using the gradient information in more sophisticated ways, but much bigger improvements are likely to result from making better use of modularity (see section 12.3).

As a biological model, back-propagation is implausible. There is no evidence that synapses can be used in the reverse direction, or that neurons can propagate error derivatives backwards (using a linear transfer function) as well as propagating activity levels forwards using a non-linear transfer function. One approach is to try to back-propagate the derivatives using separate circuitry that *learns* to have the same

¹⁰Tesauro (1987) reports a case in which the number of weight updates is roughly proportional to the number of training cases (it is actually a $4/3$ power law).

weights as the forward circuitry (Parker, 1985). A second approach, which seems to be feasible for self-supervised back-propagation, is to use a method called "recirculation" that approximates gradient descent and is much more biologically plausible (Hinton and McClelland and Goodhill, 1987). At present, back-propagation should be treated as a mechanism for demonstrating the kind of learning that can be done using gradient descent, without implying that the brain does gradient descent in the same way.

7. Boltzmann machines

A Boltzmann machine (Ackley, Hinton, and Sejnowski, 1985; Hinton and Sejnowski, 1986) is a generalization of a Hopfield net (see section 4.2) in which the units update their states according to a stochastic decision rule. The units have states of 1 or 0, and the probability that unit j adopts the state 1 is given by

$$P_j = \frac{1}{1 + e^{-\Delta E_j/T}} \quad (10)$$

where $\Delta E_j = x_j$ is the total input received by the j^{th} unit and T is the "temperature". It can be shown that if this rule is applied repeatedly to the units, the network will reach "thermal equilibrium". At thermal equilibrium the units still change state, but the *probability* of finding the network in any global state remains constant and obeys a Boltzmann distribution in which the probability ratio of any two global states depends solely on their energy difference.

$$\frac{P_A}{P_B} = e^{-(E_A - E_B)/T}$$

At high temperature, the network reaches equilibrium rapidly but low energy states are not much more probable than high energy states. At low temperature the network approaches equilibrium more slowly, but low energy states are much more probable than high energy states. The fastest way to approach low-temperature equilibrium is generally to start at a high temperature and to gradually reduce the temperature. This is called "simulated annealing" (Kirkpatrick, Gelatt and Vecchi, 1983). Simulated annealing allows Boltzmann machines to find low energy states with high probability. If some units are clamped to represent an input vector, and if the weights in the network represent the constraints of the task domain, the network can settle on a very plausible output vector given the current weights and the current input vector.

For complex tasks there is generally no way of expressing the constraints by using weights on pairwise connections between the input and output units. It is necessary to use hidden units that represent higher-order features of the domain. This creates a problem: Given a limited number of hidden

units, what higher-order features should they represent in order to approximate the required input-output mapping as closely as possible? The beauty of Boltzmann machines is that the simplicity of the Boltzmann distribution leads to a very simple learning procedure which adjusts the weights so as to use the hidden units in an optimal way.

The network is "shown" the mapping that it is required to perform by clamping an input vector on the input units and clamping the required output vector on the output units. If there are several possible output vectors for a given input vector, each of the possibilities is clamped on the output units with the appropriate probability. The network is then annealed until it approaches thermal equilibrium at a temperature of 1. It then runs for a fixed time at equilibrium and each connection measures the fraction of the time during which both the units it connects are active. This is repeated for all the various input-output pairs so that each connection can measure $\langle s_i s_j \rangle^+$, the expected probability, averaged over all cases, that unit i and unit j are simultaneously active at thermal equilibrium when the input and output vectors are both clamped.

The network must also be run in just the same way but without clamping the output units. Again, it reaches thermal equilibrium with each input vector clamped and then runs for a fixed additional time to measure $\langle s_i s_j \rangle^-$, the expected probability that both units are active at thermal equilibrium when the output vector is not clamped. Each weight is then updated by an amount proportional to the difference between these two quantities

$$\Delta w_{ij} = \varepsilon (\langle s_i s_j \rangle^+ - \langle s_i s_j \rangle^-)$$

It has been shown (Ackley, Hinton, and Sejnowski, 1985) that if ε is sufficiently small this performs gradient descent in an information theoretic measure, G , of the difference between the behavior of the output units when they are clamped and their behavior when they are not clamped.

$$G = \sum_{\alpha, \beta} P^+(I_\alpha \& O_\beta) \ln \frac{P^+(O_\beta | I_\alpha)}{P^-(O_\beta | I_\alpha)}$$

where I_α is a state vector over the input units, O_β is a state vector over the output units, P^+ is a probability measured when both the input and output units are clamped, and P^- is a probability measured at thermal equilibrium when only the input units are clamped.

G is called the *asymmetric divergence* or *Kullback information*, and its gradient has the same form for connections between input and hidden units, connections between pairs of hidden units, connections

between hidden and output units, and connections between pairs of output units. G can be viewed as the sum of two terms: The entropy of the "desired" probability distribution that is clamped on the output units and the cross-entropy between this desired probability distribution and the "actual" distribution exhibited by the output units when they are not clamped. The first term cannot be changed by altering the weights, so minimizing G is equivalent to minimizing the cross-entropy, which means that Boltzmann machines use the same cost function as one form of back-propagation (see section 6.3).

A special case of the learning procedure is when there are no input units. It can then be viewed as an unsupervised learning procedure which learns to model a probability distribution that is specified by clamping vectors on the output units with the appropriate probabilities. The advantage of modeling a distribution in this way is that the network can then perform completion. When a partial vector is clamped over a subset of the output units, the network produces completions on the remaining output units. If the network has learned the training distribution perfectly, its probability of producing each completion is guaranteed to match the environmental conditional probability of this completion given the clamped partial vector.

The learning procedure can easily be generalized to networks where each term in the energy function is the product of a weight, $w_{i,j,k,\dots}$ and an arbitrary function, $f(i,j,k,\dots)$, of the states of a subset of the units. The network must be run so that it achieves a Boltzmann distribution in the energy function. The generalized learning procedure is simply to change the weight by an amount proportional to the difference between $\langle f(i,j,k,\dots) \rangle^+$ and $\langle f(i,j,k,\dots) \rangle^-$.

The learning procedure using simple pairwise connections has been shown to produce appropriate representations in the hidden units (Ackley, Hinton and Sejnowski, 1985) and it has also been used for speech recognition (Prager, Harrison, and Fallside, 1986). However, it is considerably slower than back-propagation because of the time required to reach equilibrium in large networks. Also, the process of estimating the gradient introduces several practical problems. If the network does not reach equilibrium the estimated gradient has a systematic error, and if too few samples are taken to estimate $\langle s_i s_j \rangle^+$ and $\langle s_i s_j \rangle^-$ accurately the estimated gradient will be extremely noisy because it is the difference of two noisy estimates. Even when the noise in the estimate of the difference has zero mean, its variance is a function of $\langle s_i s_j \rangle^+$ and $\langle s_i s_j \rangle^-$. When these quantities are near zero or one, their estimates will have much lower variance than when they are near 0.5. This non-uniformity in the variance gives the hidden units a surprisingly strong tendency to develop weights that cause them to be on all the time or off all the time. A

familiar version of the same effect can be seen if sand is sprinkled on a vibrating sheet of tin. Nearly all the sand clusters at the points that vibrate the least, even though there is no bias in the direction of motion of an individual grain of sand.

One interesting feature of the Boltzmann machine is that it is relatively easy to put it directly onto a chip which has dedicated hardware for each connection and performs the annealing extremely rapidly using analog circuitry that computes the energy gap of a unit by simply allowing the incoming charge to add itself up, and makes stochastic decisions by using physical noise. Alspector and Allen (1987) have designed a chip which will run about 1 million times as fast as a simulation on a VAX, and they intend to fabricate it. Such chips may make it possible to apply connectionist learning procedures to practical problems, especially if they are used in conjunction with modular approaches that allow the learning time to scale better with the size of the task.

7.1. Maximizing reinforcement and entropy in a Boltzmann machine

The Boltzmann machine learning procedure is based on the simplicity of the expression for the derivative of the asymmetric divergence between the conditional probability distribution exhibited by the output units of a Boltzmann machine and a desired conditional probability distribution. The derivatives of certain other measures are also very simple if the network is allowed to reach thermal equilibrium. The entropy of the states of the machine is given by

$$H = - \sum_{\alpha} P_{\alpha} \ln P_{\alpha}$$

where P_{α} is the probability of a global configuration, and H is measured in units of $\log_2 e$ bits. Its derivative is

$$\frac{\partial H}{\partial w_{ij}} = \frac{1}{T} (\langle E s_i s_j \rangle - \langle E \rangle \langle s_i s_j \rangle)$$

So if each weight has access to the global energy, E , it is easy to manipulate the entropy. (11)

It is also easy to perform gradient ascent in expected reinforcement if the network is given a global reinforcement signal, R , that depends on its state. The derivative of the expected reinforcement with respect to each weight is

$$\frac{\partial R}{\partial w_{ij}} = \frac{1}{T} (\langle R s_i s_j \rangle - \langle R \rangle \langle s_i s_j \rangle)$$

A recurrent issue in reinforcement learning procedures is how to trade-off short-term optimization of(12)

expected reinforcement against the diversity required to discover actions that have a higher reinforcement than the network's current estimate. If we use entropy as a measure of diversity, and we assume that the system tries to optimize some linear combination of the expected reinforcement and the entropy of its actions, it can be shown that its optimal strategy is to pick actions according to a Boltzmann distribution, where the expected reinforcement of a state is the analog of negative energy and the parameter that determines the relative importance of expected reinforcement and diversity is the analog of temperature. This result follows from the fact that the Boltzmann distribution is the one which maximizes entropy for a given expected energy.

This suggests a learning procedure in which the system represents the expected value of an action by its negative energy, and picks actions by allowing a Boltzmann machine to reach thermal equilibrium. If the weights are updated using equations Eq. 11 and Eq. 12 the negative energies of states will tend to become proportional to their expected reinforcements, since this is the way to make the derivative of H balance the derivative of R . Once the system has learned to represent the reinforcements correctly, variations in the temperature can be used to make it more or less conservative in its choice of actions whilst always making the optimal trade-off between diversity and expected reinforcement. Unfortunately, this learning procedure does not make use of the most important property of Boltzmann machines which is their ability to compute the quantity $\langle s_i s_j \rangle$ given some specified state of the output units. Also, it is much harder to compute the derivative of the entropy if we are only interested in the entropy of the state vectors over the output units.

8. Maximizing mutual information: A semi-supervised learning procedure

Supervised learning procedures usually confound two different kinds of information: Information about which of several possible distributions the current input vector came from (i.e. its category), and information about the "desired" states that certain units ought to adopt (i.e. the desired representation of that category in the output units). One "semi-supervised" method of training a hidden unit is to provide it with information about what category the input vector came from, but to refrain from specifying the state that it ought to adopt. Instead, its incoming weights are modified so as to maximize the information that the state of the hidden unit provides about the category of the input vector. The derivative of this mutual information is relatively easy to compute and so it can be maximized by gradient ascent (Pearlmutter and Hinton, 1986). For difficult discriminations that cannot be performed in a single step this is a good heuristic for producing encodings of the input vector that allow the discrimination to be made more easily.

Figure 8 shows an example of a difficult two-way discrimination and illustrates the kinds of discriminant function that maximize the information provided by the state of the hidden unit.

.....
Insert figure 8 about here
.....

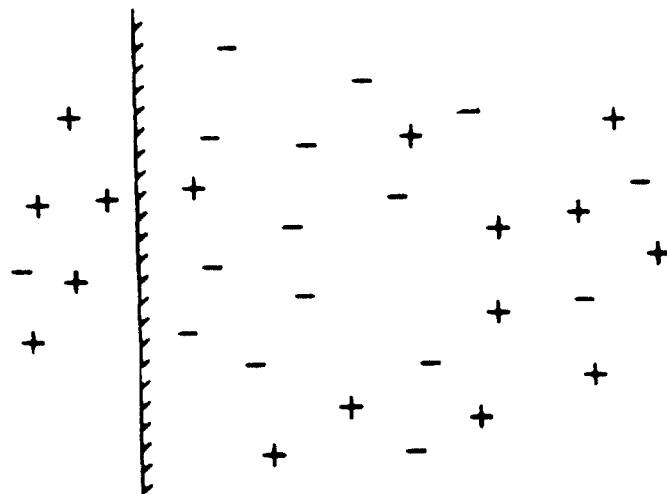
If each hidden unit within a layer independently maximizes the mutual information between its state and the category of the input vector, many units are likely to discover similar, highly-correlated features. One way to force the hidden units to diversify is to make each hidden unit receive its inputs from a different subset of the units in the layer below. A second method is to ignore cases in which the input vector is correctly classified by the output units and to maximize the mutual information between the state of the hidden unit and the category of the input *given that the input is incorrectly classified*.¹¹

If the two input distributions that must be discriminated consist of examples taken from some structured domain and examples generated at random (but with the same first-order statistics as the structured domain), this semi-supervised procedure will discover higher-order features that characterize the structured domain and so it can be made to act like the type of unsupervised learning procedure described in section 9.

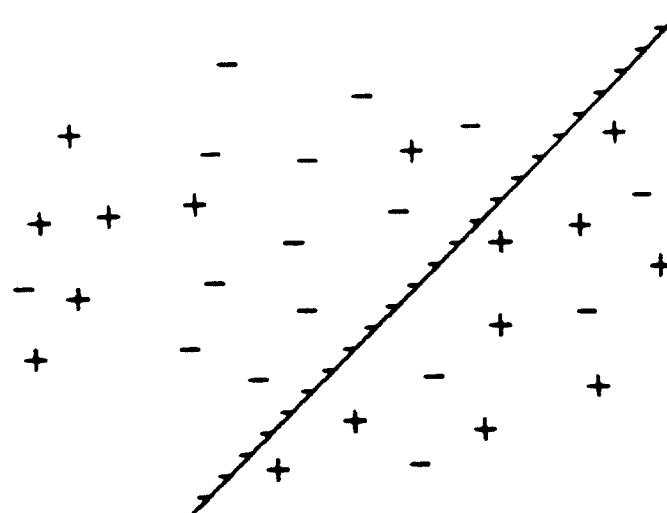
9. Unsupervised Hebbian learning

The central idea of unsupervised Hebbian learning is that an individual unit can develop selectivity to an "important" subset of its ensemble of input vectors by using a simple weight modification procedure that depends on the correlation of pre-synaptic and post-synaptic activity. Typical examples of this kind of learning are described by Cooper, Liberman and Oja (1979) and by Bienenstock, Cooper, and Munro (1982). A criticism of early versions of this approach, from a computational point of view, was that the researchers often postulated a simple synaptic modification rule and then explored its consequences rather than rigorously specifying the computational goal and then deriving the appropriate synaptic modification rule. However, an important recent development unifies these two approaches by showing that a relatively simple Hebbian rule can be viewed as the gradient of an interesting function. The function can therefore be viewed as a specification of what the learning is trying to achieve.

¹¹This method of weighting the statistics by some measure of the overall error or importance of a case can often be used to allow global measures of the performance of the whole network to influence local, unsupervised learning procedures.



(a)



(b)

Figure 8: (a) There is high mutual information between the state of a binary threshold unit that uses the hyperplane shown and the distribution (+ or -) that the input vector came from. (b) The probability, given that the unit is on, that the input came from the "+" distribution is not as high using the diagonal hyperplane. However, the unit is on more often. Other things being equal, a unit conveys most mutual information if it is on half the time.

9.1. A recent development of unsupervised Hebbian learning

In a recent series of papers Linsker has shown that with proper normalization of the weight changes, an unsupervised Hebbian learning procedure in which the weight change depends on the correlation of pre-synaptic and post-synaptic activity can produce a surprising number of the known properties of the receptive fields of neurons in visual cortex, including center-surround fields (Linsker, 1986a) orientation-tuned fields (Linsker, 1986b) and orientation columns (Linsker, 1986c). The procedure operates in a multilayer network in which there is innate spatial structure so that the inputs to a unit in one layer tend to come from nearby locations in the layer below. Linsker demonstrates that the emergence of biologically suggestive receptive fields depends on the relative values of a few generic parameters. He also shows that for each unit, the learning procedure is performing gradient ascent in a measure whose main term is the ensemble average (across all the various patterns of activity in the layer below) of:

$$\sum_{ij} w_i s_i w_j s_j$$

where w_i and w_j are the weights on the i^{th} and j^{th} input lines of a unit and s_i and s_j are the activities on those input lines.

It is not initially obvious why maximizing the pairwise covariances of the weighted activities produces receptive fields that are useful for visual information processing. Linsker does not discuss this question in his original three papers. However, he has now shown (Linsker, 1987) that the learning procedure maximizes the variance in the activity of the post-synaptic unit subject to a "resource" constraint on overall synaptic strength. Since Linsker's units are linear, maximizing the post-synaptic variance is equivalent to extracting the first component in a principal components analysis. This component is the one that would minimize the sum-squared reconstruction error if we tried to reconstruct the activity vector of the presynaptic units from the activity level of the post-synaptic unit. Thus we can view Linsker's learning procedure as a way of ensuring that the activity of a unit conveys as much information as possible about its pre-synaptic input vector. A similar analysis can be applied to competitive learning (see section 10).

10. Competitive learning

Competitive learning is an unsupervised procedure that divides a set of input vectors into a number of disjoint clusters in such a way that the input vectors within each cluster are all similar to one another. It is called competitive learning because there is a set of hidden units which compete with one another to

become active. There are many variations of the same basic idea, and only the simplest version is described here. When an input vector is presented to the network, the hidden unit which receives the greatest total input wins the competition and turns on with an activity level of 1. All the other hidden units turn off. The winning unit then adds a small fraction of the current input vector to its weight vector. So, in future, it will receive even more total input from this input vector. To prevent the same hidden unit from being the most active in all cases, it is necessary to impose a constraint on each weight vector that keeps the sum of the weights (or the sum of their squares) constant. So when a hidden unit becomes more sensitive to one input vector it becomes less sensitive to other input vectors.

Rumelhart and Zipser (1985) present a simple geometrical model of competitive learning. If each input vector has three components and is of unit length it can be represented by a point on the surface of the unit sphere. If the weight vectors of the hidden units are also constrained to be of unit length, they too can be represented by points on the unit sphere as shown in figure 9. The learning procedure is equivalent to finding the weight vector that is closest to the current input vector, and moving it closer still by an amount that is proportional to the distance. If the weight changes are sufficiently small, this process will stabilize when each weight vector is at the center of gravity of a cluster of input vectors.

Insert figure 9 about here

We can think of the network as performing the following task: Represent the current input vector, s_c , as accurately as possible by using a single active hidden unit. The representation is simply the weight vector, w_c , of the hidden unit which is active in case c . If the weight changes are sufficiently small, this version of competitive learning performs steepest descent in a measure of the mean squared inaccuracy of the representation. The solutions it finds are minima of the function

$$E = \frac{1}{2} \sum_c (w_c - s_c)^2$$

Although they use the geometrical analogy described above, Rumelhart and Zipser actually use a slightly different learning rule which cannot be interpreted as performing steepest descent in such a simple error function.

There are many variations of competitive learning in the literature (Von der Malsburg, 1973; Fukushima, 1975; Grossberg, 1976) and there is not space here to review them all. One major theme has been to show that competitive learning can produce topographic maps if more than one hidden unit is

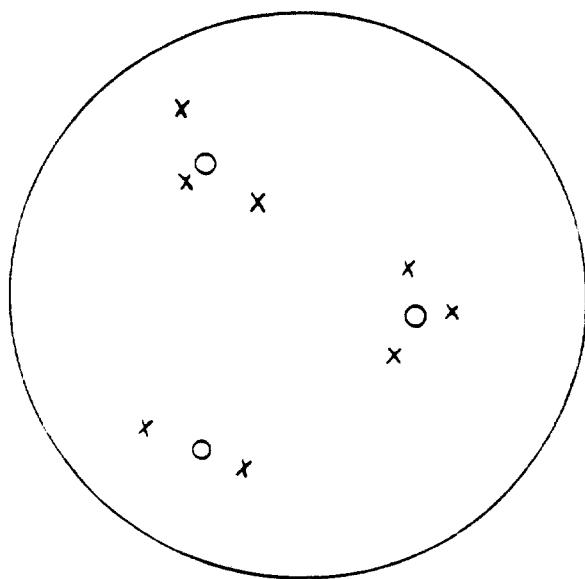


Figure 9: The input vectors are represented by points marked "x" on the surface of a sphere. The weight vectors of the hidden units are represented by points marked "o". After competitive learning, each weight vector will be close to the center of gravity of a cluster of input vectors.

allowed to be active at a time and physically adjacent pairs of hidden units inhibit each other less than more distant pairs (Kohonen, 1982; Amari, 1983). A model with similarities to competitive learning has been used by Willshaw and Von der Malsburg (1979) to explain the formation of topographic maps between the retina and the tectum. Recently, it has been shown that a variation of this model can be interpreted as performing steepest descent in an error function and can be applied to a range of optimization problems that involve topographic mappings between geometrical structures (Durbin and Willshaw, 1987).

A related theme is to show that competitive learning can perform dimensionality reduction, so that surplus degrees of freedom are removed from the input vector and it is represented accurately by a point in a lower dimensional space (Kohonen, 1982). It is not clear how this compares in efficiency with self-supervised back-propagation (see section 6.5) for dimensionality reduction.

Fukushima and Miyake (1982) have demonstrated that a version of competitive learning can be used to allow a multilayer network to recognize simple two-dimensional shapes in a number of different positions. After learning, the network can recognize a familiar shape in a novel position. The ability to generalize across position depends on using a network in which the layers of units that learn are interleaved with layers of non-learning units which are pre-wired to generalize across position. Thus, the network does not truly learn translation invariance. By contrast, it is possible to design a back-propagation network that starts with no knowledge of the effects of translation and no knowledge of which input units are adjacent in the image. After sufficient experience, the network can correctly identify familiar, simple shapes in novel positions (Hinton, 1987c).

10.1. The relationship between competitive learning and back-propagation

Because it is performing gradient descent in a measure of how accurately the input vector could be reconstructed, competitive learning has a close relationship to self-supervised back-propagation. Consider a three layer encoder network in which the desired states of the output units are the same as the actual states of the input units. Suppose that each weight from an input unit to a hidden unit is constrained to be identical to the weight from that hidden unit to the corresponding output unit. Suppose, also, that the output units are linear and the hidden units, instead of using the usual non-linear transfer function, use the same "winner-take-all" non-linearity as is used in competitive learning. So only one hidden unit will be active at a time, and the actual states of the output units will equal the weights of the

active hidden unit. This makes it is easy to compute the error derivatives of the weights from the hidden units to the output units. For weights from the active hidden unit the derivatives are simply proportional to the difference between the actual and desired outputs (which equals the difference between the weight and the corresponding component of the input vector). For weights from inactive hidden units the error derivatives are all zero.

Normally, back-propagation is needed in order to compute the error derivatives of the weights from the input units to the hidden units, but the winner-take-all non-linearity makes back-propagation unnecessary in this network because all these derivatives are equal to zero. So long as the same hidden unit wins the competition, its activity level is not changed by changing its input weights. At the point where a small change in the weights would change the winner from one hidden unit to another, both hidden units fit the input vector equally well, so changing winners does not alter the total error in the output (even though it may change the output vector a lot). Because the error derivatives are so simple, we can still do the learning if we omit the output units altogether. This removes the output weights, and so we no longer need to constrain the input and output weights of a hidden unit to be identical. Thus the simplified version of competitive learning is a degenerate case of self-supervised back-propagation.

It would be interesting if a mechanism as simple as competitive learning could be used to implement gradient descent in networks that allow the m most activated hidden units to become fully active (where $m > 1$). This would allow the network to create more complex, distributed representations of the input vectors. Unfortunately the implementation is not nearly as simple because it is no longer possible to omit the output layer. The output units are needed to combine the effects of all the active hidden units and compare the combined effect with the input vector in order to compute the error derivatives of the output weights. Also, at the point at which one hidden unit ceases to be active and another becomes active, there may be a large change in the total error, so at this point there is a discontinuity in the error derivatives of the weights from the input to the hidden units. It thus appears that the simplicity of the mechanism required for competitive learning is crucially dependent on the fact that only one hidden unit within a group is active.

11. Reinforcement learning procedures

There is a large and complex literature on reinforcement learning procedures which is beyond the scope of this paper. The main aim of this section is to give an informal description of a few of the recent ideas in the field that reveals their relationship to other types of connectionist learning.

A central idea in many reinforcement learning procedures is that we can assign credit to a local decision by *measuring* how it correlates with the global reinforcement signal. Various different values are tried for each local variable (such as a weight or a state), and these variations are correlated with variations in the global reinforcement signal. Normally, the local variations are the result of independent stochastic processes, so if enough samples are taken each local variable can average away the noise caused by the variation in the other variables to reveal its own effect on the global reinforcement signal (given the current expected behavior of the other variables). The network can then perform gradient ascent in the expected reinforcement by altering the probability distribution associated with each local variable in the direction that increases the expected reinforcement. If the probability distributions are altered after each trial, the network performs a stochastic version of gradient ascent.

The main advantage of reinforcement learning is that it is easy to implement because, unlike back-propagation which *computes* the effect of changing a local variable, the "credit assignment" does not require any special apparatus for communicating between the local variables. So reinforcement learning can be used in complex systems which would be very hard to invert analytically. The main disadvantage is that it is very inefficient when there are more than a few local variables. Even in the trivial case when all the local variables contribute independently to the global reinforcement signal, $\text{order}(NM)$ trials are required to allow all of the M values of a variable to achieve a reasonable signal-to-noise ratio by averaging away the noise caused by the N other variables. So reinforcement learning is very inefficient for large systems unless they are divided into smaller modules. It is as if each person in the United States tried to decide whether he or she had done a useful day's work by observing the gross national product on a day by day basis.

A second disadvantage is that gradient ascent may get stuck in local optima. As a network concentrates more and more of its trials on combinations of values that give the highest expected reinforcement, it gets less and less information about the reinforcements caused by other combinations of values.

11.1. Delayed reinforcement

In many real systems, there is a delay between an action and the resultant reinforcement, so in addition to the normal problem of deciding how to assign credit to decisions about hidden variables, there is a temporal credit assignment problem (Sutton, 1984). In the iterative version of back-propagation (section 6.4), temporal credit assignment is performed by explicitly computing the effect of each activity level on the eventual outcome. In reinforcement learning procedures, temporal credit assignment is typically performed by learning to associate "secondary" reinforcement values with the states that are intermediate in time between the action and the external reinforcement. One important idea is to make the reinforcement value of an intermediate state regress towards the weighted average of the reinforcement values of its successors, where the weightings reflect the conditional probabilities of the successors. In the limit, this causes the reinforcement value of each state to be equal to the expected reinforcement of its successor, and hence equal to the expected final reinforcement.¹² Sutton (1987) explains why, in a stochastic system, it is typically more efficient to regress towards the reinforcement value of the next state rather than the reinforcement value of the final outcome. Barto, Sutton and Anderson (1983) have demonstrated the usefulness of this type of procedure for learning with delayed reinforcement.

11.2. The A_{R-P} procedure

One obvious way of mapping results from learning automata theory onto connectionist networks is to treat each unit as an automaton and to treat the states it adopts as its actions. Barto and Anandan (1985) describe a learning procedure of this kind called "Associative Reward-Penalty" or A_{R-P} , which uses stochastic units like those in a Boltzmann machine (see equation 10). They prove that if the input vectors are linearly independent and the network only contains one unit, A_{R-P} finds the optimal values of the weights. They also show empirically that if the same procedure is applied in a network of such units, the hidden units develop useful representations. Williams (1986) has shown that a limiting case of the A_{R-P} procedure performs stochastic gradient ascent in expected reinforcement.

¹²There may also be a "tax" imposed for failing to achieve the external reinforcement quickly. This can be implemented by reducing the reinforcement each time it is regressed to an earlier state.

11.3. Achieving global optimality by reinforcement learning

Thatachar and Sastry (1985) use a different mapping between automata and connectionist networks. Each *connection* is treated as an automaton and the weight-values that it takes on are its actions. On each trial, each connection chooses a weight (from a discrete set of alternatives) and then the network maps an input vector into an output vector and receives positive reinforcement if the output is correct. They present a learning procedure for updating the probabilities of choosing particular weight-values. If the probabilities are changed slowly enough, the procedure is guaranteed to converge on the globally optimal combination of weights, even if the network has hidden layers. Unfortunately their procedure requires exponential space because it involves storing and updating a table of estimated expected reinforcements that contains one entry for every combination of weights.

11.4. The relative payoff procedure

If we are content to reach a local optimum, it is possible to use a very simple learning procedure. I shall describe a connectionist implementation that uses yet another way of mapping automata onto connectionist networks. Each connection is treated as a stochastic switch that has a certain probability of being closed at any moment (Minsky, 1954). If the switch is open, the "post-synaptic" unit receives an input of 0 along that connection, but if the switch is closed it transmits the state of the "pre-synaptic" unit. A synapse can be modeled as a set of these stochastic switches arranged in parallel. Each unit computes some fixed function of the vector of inputs that it receives on its incoming connections. Learning involves altering the switch probabilities to maximize the expected reinforcement signal.

A learning procedure called L_{R_i} can be applied in such networks. It is only guaranteed to find a local optimum of the expected reinforcement, but it is very simple and does not require any communication between the switches. A "trial" consists of four stages:

1. Set the switch configuration. For each switch in the network, decide whether it is open or closed on this trial using the current switch probability. The decisions are made independently for all the switches.
2. Run the network with this switch configuration. There are no constraints on the connectivity so cycles are allowed, and the units can also receive external inputs at any time. The constraint on the external inputs is that the probability distribution over patterns of external input must be independent of the switch settings and it must be stationary.

3. Compute the reinforcement signal. This can be any non-negative, stationary function of the behavior of the network and of the external input it received during the trial.

4. Update the switch probabilities. For each switch that was closed during the trial, we increment its probability by $\epsilon R(1-p)$, where R is the reinforcement produced by the trial, p is the switch probability and ϵ is a small coefficient. For each switch that was open, we decrement its probability by ϵRp .

If ϵ is sufficiently small this procedure is guaranteed to hill-climb in expected reinforcement. The "batch" version of the procedure involves observing the reinforcement signal over a large number of trials before updating the switch probabilities. If a sufficient number of trials are observed, the following "relative payoff" update procedure always increases expected reinforcement: Change the switch probability to be equal to the fraction of the total reinforcement received when the switch was closed. This can cause large changes in the probabilities, and I know of no proof that it hill-climbs in expected reinforcement, but in practice it always works. The direction of the jump in the switch-probability-space caused by the batch version of the procedure is the same as the expected direction of the small change in switch probabilities caused by the "on-line" version.

The "on-line" version of this procedure has been proposed as a model of perceptual interpretation (Harth, Pandya and Unnikrishnan, 1986). Williams (1986b) gives a detailed analysis of a similar learning procedure.

A variation of the relative payoff procedure can be used if the goal is to make the "responses" of a network match some desired probability distribution rather than maximize expected reinforcement. We simply define the reinforcement signal to be the desired probability of a response divided by the network's current probability of producing that response. If a sufficient number of trials are made before updating the switch probabilities, it can be shown (Larry Gillick and Jim Baker, personal communication) that this procedure is guaranteed to decrease an information theoretic measure of the difference between the desired probability distribution over responses and the actual probability distribution. The measure is actually the G measure described in the section on Boltzmann machines and the proof is an adaptation of the proof of the EM procedure (Dempster, Laird and Rubin, 1976).

11.5. Genetic algorithms

Holland and his co-workers (Holland, 1975; Grefenstette, 1985) have investigated a class of learning procedures which they call "genetic algorithms" because they are explicitly inspired by an analogy with evolution. Genetic algorithms operate on a population of individuals to produce a better adapted population. In the simplest case, each individual member of the population is a binary vector, and the two possible values of each component are analogous to two alternative versions (alleles) of a gene. There is a fitness function which assigns a real-valued fitness to each individual and the aim of the "learning" is to raise the average fitness of the population. New individuals are produced by choosing two existing individuals as parents (with a bias towards individuals of higher than average fitness) and copying some component values from one parent and some from the other. Holland (1975) has shown that for a large class of fitness functions, this is an effective way of discovering individuals that have high fitness.

11.6. Genetic learning and the relative payoff rule

If an entire generation of individuals is simultaneously replaced by a generation of their offspring, genetic learning has a close relationship to the batch form of the $L_{R,I}$ procedure described in section 11.4. This is most easily understood by starting with a particularly simple version of genetic learning in which every individual in generation $t+1$ has many different parents in generation t . Candidate individuals for generation $t+1$ are generated from the existing individuals in generation t in the following way: To decide the value of the i^{th} component of a candidate, we randomly choose one of the individuals in generation t and copy the value of its i^{th} component. So the probability that the i^{th} component of a candidate has a particular value is simply the relative frequency of that value in generation t . A selection process then operates on the candidates: Some are kept to form generation $t+1$ and others are discarded. The fitness of a candidate is simply the probability that it is not discarded by the selection process. Candidates that are kept can be considered to have received a reinforcement of 1 and candidates that are discarded receive a reinforcement of 0. After selection, the probability that the i^{th} component has a particular value is equal to the fraction of the successful candidates that have that value. This is exactly the relative payoff rule described in section 11.4. The probabilities it operates on are the relative frequencies of alleles in the population instead of switch probabilities.

If the value of every component is determined by an independently chosen parent, all information about the correlations between the values of different components is lost when generation $t+1$ is produced from generation t . If, however, we use just two parents we maximize the tendency for the

pairwise and higher-order correlations to be preserved. This tendency is further increased if components whose correlations are important are near one another and the values of nearby components are normally taken from the same parent. So a population of individuals can effectively represent the probabilities of small combinations of component values as well as the probabilities of individual values. Genetic learning works well when the fitness of an individual is determined by these small combinations, which Holland calls critical schemas.

11.7. Iterated genetic hill-climbing

It is possible to combine genetic learning with gradient descent (or hill-climbing) to get a hybrid learning procedure called "iterated genetic hill-climbing" or "IGH" that works better than either learning procedure alone (Brady, 1985; Ackley, 1987). IGH is as a form of multiple-restart hill-climbing in which the starting points, instead of being chosen at random, are chosen by "mating" previously discovered local optima. Alternatively, it can be viewed as genetic learning in which each new individual is allowed to perform hill-climbing in the fitness function before being evaluated and added to the population. Ackley (1987) shows that a stochastic variation of IGH can be implemented in a connectionist network that is trying to learn which output vector produces a high enough payoff to satisfy some external criterion.

12. Discussion

This review has focused on a small number of recent connectionist learning procedures. There are many other interesting procedures which have been omitted (Grossberg 1980; Volper and Hampson, 1986; Hampson and Volper, 1987; Hogg and Huberman, 1984; Kerszberg and Bergman, 1986; Edelman and Reeke, 1982). Rather than attempting to cover all of these I conclude by discussing two major problems that plague most of the procedures I have described.

12.1. Generalization

A major goal of connectionist learning is to produce networks that generalize correctly to new cases after training on a sufficiently large set of typical cases from some domain. In much of the research, there is no formal definition of what it means to generalize correctly. The network is trained on examples from a domain that the experimenter understands (like the family relationships domain described in section 6) and it is judged to generalize correctly if its generalizations agree with those of the experimenter. This is sufficient as an informal demonstration that the network can indeed perform non-trivial generalization, but it gives little insight into the reasons why the generalizations of the network and the experimenter agree.

and so it does not allow predictions to be made about when networks will generalize correctly and when they will fail.

What is needed is a formal theory of what it means to generalize correctly. One approach that has been used in studying the induction of grammars is to define a hypothesis space of possible grammars, and to show that with enough training cases the system will converge on the correct grammar with probability 1. This is called "identification in the limit" (Angluin and Smith, 1984).

Valiant (1984) has recently introduced a rather more subtle criterion of success in order to distinguish classes of boolean function that can be induced from examples in polynomial time from classes that require exponential time. He assumes that the hypothesis space is known in advance and he allows the training cases to be selected according to any stationary distribution but insists that the same distribution be used to generate the test cases. The induced function is considered to be good enough if it differs from the true function on less than a small fraction, $1/h$, of the test cases. A class of boolean functions is considered to be learnable in polynomial time if, for any choice of h , there is a probability of at least $(1 - 1/h)$ that the induced function is good enough after a number of training examples that is polynomial in both h and the number of arguments of the boolean function. Using this definition, Valiant has succeeded in showing that several interesting subclasses of boolean function are learnable in polynomial time. Our understanding of other connectionist learning procedures would be considerably improved if we could derive similar results that were as robust against variations in the distribution of the training examples.

The work on inducing grammars or boolean functions may not provide an appropriate framework for studying systems that learn inherently stochastic functions, but the general idea of starting with a hypothesis space of possible functions carries over. A widely used statistical approach involves maximizing the *a posteriori* likelihood of the model (i.e. the function) given the data. If the data really is generated by a function in the hypothesis space and if the amount of information in the training data greatly exceeds the amount of information required to specify a point in the hypothesis space, the maximum likelihood function is very probably the correct one, so the network will then generalize correctly. Some connectionist learning schemes (e.g. the Boltzmann machine learning procedure) can be made to fit this approach exactly. If a Boltzmann machine is trained with much more data than there are weights in the machine, and if it really does find the global minimum of G , and if the correct answer lies in

the hypothesis space (which is defined by the architecture of the machine)¹³, then there is every reason to suppose that it will generalize correctly, even if it has only been trained on a small fraction of the possible cases. Unfortunately, this kind of guarantee is of little use for practical problems where we usually know in advance that the "true" model does not lie in the hypothesis space of the network. What needs to be shown is that the best available point within the hypothesis space (even though its not a perfect model) will also generalize well to test cases. One approach to this issue is described by Natarajan (1987).

A simple thought experiment shows that the "correct" generalization from a set of training cases, however it is defined, must depend on how the input and output vectors are encoded. Consider a mapping, M_I , from entire input vectors onto entire input vectors and a mapping, M_O , from entire output vectors onto entire output vectors. If we introduce a pre-coding stage that uses M_I and a post-coding stage that uses M_O we can convert a network that generalizes in one way into a network that generalizes in any other way we choose simply by choosing M_I and M_O appropriately.

12.2. The speed of learning

Most existing connectionist learning procedures are slow, particularly procedures that construct complicated internal representations. One way to speed them up is to use optimization methods such as recursive least squares that converge faster. If the second derivatives can be computed or estimated they can be used to pick a direction for the weight change vector that yields faster convergence than the direction of steepest descent (Parker, 1987). It remains to be seen how well such methods work for the error surfaces generated by multilayer networks learning complex tasks.

A second method of speeding up learning is to use dedicated hardware for each connection and to map the inner-loop operations into analog instead of digital hardware. As Alspector and Allen (1987) have demonstrated, the speed of one particular learning procedure can be increased by a factor of about a million if we combine these techniques. This significantly increases our ability to explore the behavior of relatively small systems, but it is not a panacea. By using silicon in a different way we typically gain a large but constant factor (optical techniques may eventually yield a huge constant factor), and by

¹³One popular idea is that evolution implicitly chooses an appropriate hypothesis space by constraining the architecture of the network and learning then identifies the most likely hypothesis within this space. How evolution arrives at sensible hypothesis spaces in reasonable time is usually unspecified. The evolutionary search for good architectures may actually be guided by learning (Hinton and Nowlan, 1987).

dedicating a processor to each of the N connections we gain at most a factor of N in time at the cost of at least a factor of N in space. For a learning procedure with a time complexity of, say, $O(N \log N)$ a speed-up of N makes a very big difference. For a procedure with a complexity of, say, $O(N^3)$ alternative technologies and parallelism will help significantly for small systems, but not for large ones.¹⁴

12.3. Hardware modularity

One of the best and commonest ways of fighting complexity is to introduce a modular, hierarchical structure in which different modules are only loosely coupled (Simon, 1969). Pearl (1986) has shown that if the interactions between a set of probabilistic variables are constrained to form a tree structure, there are efficient parallel methods for estimating the interactions between "hidden" variables. The leaves of the tree are the observables and the higher level nodes are hidden. The probability distribution for each variable is constrained by the values of its immediate parents in the tree. Pearl shows that these conditional probabilities can be recovered in time $O(N \log N)$ from the pairwise correlations between the values of the leaves of the tree. Remarkably, it is also possible to recover the tree structure itself in the same time.

Self-supervised back-propagation (see section 6.5) was originally designed to allow efficient bottom-up learning in domains where there is hierarchical modular structure. Consider, for example, an ensemble of input vectors that are generated in the following modular way: Each module has a few high-level variables whose values help to constrain the values of a larger number of low-level variables. The low-level variables of each module are partitioned into several sets, and each set is identified with the high-level variables of a lower module as shown in figure 10.

Insert figure 10 about here

Now suppose that we treat the values of all the low-level variables of the leaf modules as a single input vector. Given a sufficiently large ensemble of input vectors and an "innate" knowledge of the architecture of the generator, it should be possible to recover the underlying structure by using self-supervised back-propagation to learn compact codes for the low-level variables of each leaf module. It is possible to learn codes for all the lowest level modules in parallel. Once this has been done, the network can learn codes at the next level up the hierarchy. The time taken to learn the whole hierarchical

¹⁴Tsotsos (1987) makes similar arguments in a discussion of the space-complexity of vision.

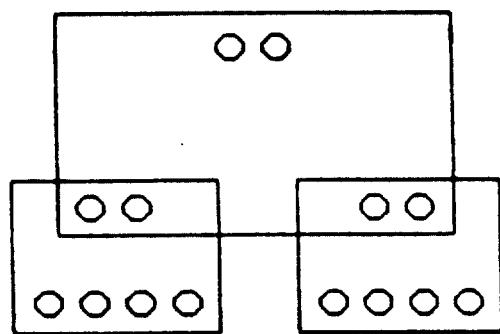


Figure 10: The lower-level variables of a high level module are the higher-level variables of several low level modules.

structure (given parallel hardware) is just proportional to the depth of the tree and hence it is $O(\log N)$ where N is the size of the input vector. An improvement on this strictly bottom-up scheme is described by Ballard (1987). He shows why it is helpful to allow top-down influences from more abstract representations to less abstract ones, and presents a working simulation.

12.4. Other types of modularity

There are several other helpful types of modularity that do not necessarily map so directly onto modular hardware but are nevertheless important for fast learning and good generalization. Consider a system which solves hard problems by creating its own subgoals. Once a subgoal has been created, the system can learn how best to satisfy it and this learning can be useful (on other occasions) even if it was a mistake to create that subgoal on this particular occasion. So the assignment of credit to the decision to create a subgoal can be decoupled from the assignment of credit to the actions taken to achieve the subgoal. Since the ability to achieve the subgoals can be learned separately from the knowledge about when they are appropriate, a system can use achievable subgoals as building blocks for more complex procedures. This avoids the problem of learning the complex procedures from scratch. It may also constrain the way in which the complex procedures will be generalized to new cases, because the knowledge about how to achieve each subgoal may already include knowledge about how to cope with variations. By using subgoals we can increase modularity and improve generalization even in systems which use the very same hardware for solving the subgoal as was used for solving the higher level goal.

There is another type of relationship between easy and hard tasks that can facilitate learning. Sometimes a hard task can be decomposed into a set of easier constituents, but other times a hard task may just be a version of an easier task that requires finer discrimination. For example, throwing a ball in the general direction of another person is much easier than throwing it through a hoop, and a good way to train a system to throw it through a hoop is to start by training it to throw it in the right general direction. This relation between easy and hard tasks is used extensively in "shaping" the behavior of animals and should also be useful for connectionist networks (particularly those that use reinforcement learning). It resembles the use of multi-resolution techniques to speed up search in computer vision (Terzopoulos, 1984). Having learned the coarse task, the weights should be close to a point in weight-space where minor adjustments can tune them to perform the finer task.

One application where this technique should be helpful is in learning filters that discriminate

between very similar sounds. The approximate shapes of the filters can be learned using spectrograms that have low resolution in time and frequency, and then the resolution can be increased to allow the filters to resolve fine details. By introducing a "regularization" term that penalizes filters which have very different weights for adjacent cells in the high-resolution spectrogram, it may be possible to allow filters to "attend" to fine detail when necessary without incurring the cost of estimating all the weights from scratch. The regularization term encodes prior knowledge that good filters should generally be smooth and so it reduces the amount of information that must be extracted from the training data (Lang and Hinton, in preparation).

12.5. Conclusion

There are now many different connectionist learning procedures that can construct appropriate internal representations in small domains, and it is likely that many more variations will be discovered in the next few years. Major new advances can be expected on a number of fronts: Techniques for making the learning time scale better may be developed; attempts to apply connectionist procedures to difficult tasks like speech recognition may actually succeed; new technologies may make it possible to simulate much larger networks; and finally the computational insights gained from studying connectionist systems may prove useful in interpreting the behavior of real neural networks.

Acknowledgements

I thank David Ackley, Joshua Alspicer, Jim Baker, Dana Ballard, Andrew Barto, John Bridle, Yann Le Cun, Mark Derthick, Larry Gillick, Ralph Linsker, Steve Nowlan, David Plaut, David Rumelhart, Terry Sejnowski, John Tsotsos and members of the Carnegie-Mellon Boltzmann group for helpful comments.

References

- Ackley, D. H., Hinton, G. E., Sejnowski, T. J. (1985). A learning algorithm for Boltzmann machines. *Cognitive Science*, 9, 147-169.
- Ackley, D. H. (1987). *Stochastic Iterated Genetic Hill-climbing* PhD Thesis). Carnegie-Mellon University, Pittsburgh PA.
- Alspector, J. & Allen, R. B. (1987). A neuromorphic VLSI learning system. In P. Loseleben (Ed.), *Advanced Research in VLSI: Proceedings of the 1987 Stanford Conference..* Cambridge, Mass.: MIT Press.
- Amari, S-I. (1967). A theory of adaptive pattern classifiers. *IEEE Transactions on Electronic Computers*, EC-16, 299-307.
- Amari, S-I. (1983). Field theory of self-organizing neural nets. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13, 741-748.
- Anderson, J. A., & Hinton, G. E. (1981). Models of information processing in the brain. In G. E. Hinton & J. A. Anderson (Eds.), *Parallel models of associative memory*. Hillsdale, NJ: Erlbaum.
- Anderson, J. A. & Mozer, M. C. (1981). Categorization and selective neurons. In G. E. Hinton & J. A. Anderson (Eds.), *Parallel Models of Associative Memory*. Hillsdale, NJ: Erlbaum.
- Bahl, L. R., Jelinek, F., & Mercer, R. L. (1983). A maximum likelihood approach to continuous speech recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-5, 179-190.
- Ballard, D. H., Hinton, G. E., & Sejnowski, T. J. (1983). Parallel visual computation. *Nature*, 306, 21-26.
- Ballard, D. H. (1986). Cortical connections and parallel processing: Structure and function. *The Behavioral and Brain Sciences*, 9, 67-120.
- Ballard, D. H. (1987). Modular learning in neural networks. *Proc. AAAI-87*. Seattle, WA.
- Barlow, H. B. (1972). Single units and sensation: A neuron doctrine for perceptual psychology? *Perception*, 1, 371-394.
- Barto, A. G. Sutton, R. S. & Anderson, C. W. (1983). Neuronlike elements that solve difficult learning control problems. *IEEE Transactions on Systems, Man, & Cybernetics*.
- Bienenstock, E. L., Cooper, L. N., & Munro, P. W. (1982). Theory for the development of neuron selectivity: Orientation specificity and binocular interaction in visual cortex. *Journal of Neuroscience*, 2, 32-48.
- Brady, R. M. (1985). Optimization strategies gleaned from biological evolution. *Nature*, 317, 804-806.

- Crick, F. & Mitchison, G. (1983). The function of dream sleep. *Nature*, 304, 111-114.
- Dempster, A.P., Laird, N. M., & Rubin, D. B. (1976). Maximum likelihood from incomplete data via the EM algorithm. *Proceedings of the Royal Statistical Society*, pp. 1-38.
- Edelman, G. M. & Reeke, G. N. (1982). Selective networks capable of representative transformations, limited generalizations, and associative memory. *Proceedings of the National Academy of Sciences USA*, 79, 2091-2095.
- Elman, J. L. and Zipser, D. (1987). *Discovering the hidden structure of speech* (Tech. Rep.). Institute for Cognitive Science Technical Report No. 8701. University of California, San Diego.,
- Feldman, J. A. & Ballard, D. H. (1982). Connectionist models and their properties. *Cognitive Science*, 6, 205-254.
- Feldman, J. A. (1986). *Neural representation of conceptual knowledge* Technical Report TR189). Department of Computer Science, University of Rochester, Rochester NY,
- Fukushima, K. (1975). Cognitron: A self-organizing multilayered neural network. *Biological Cybernetics*, 20, 121-136.
- Fukushima, K., & Miyake, S. (1982). Neocognitron: A new algorithm for pattern recognition tolerant of deformations and shifts in position. *Pattern Recognition*, 15, 455-469.
- Geman, S., & Geman D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6, 721-741.
- Grefenstette, J. J. (1985). *Proceedings of an international conference on genetic algorithms and their applications*.
- Grossberg, S. (1976). Adaptive pattern classification and universal recoding, I: Parallel development and coding of neural feature detectors. *Biological Cybernetics*, 23, 121-134.
- Grossberg, S. (1980). How does the brain build a cognitive code? *Psychological Review*, 87, 1-51.
- Hampson, S. E. & Volper, D. J. (1987). Disjunctive models of Boolean category learning. *Biological Cybernetics*, 55, 1-17.
- Harth, E., Pandya, A. S. & Unnikrishnan, K. P. (1986). Perception as an optimization process. *Proc. IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*. Miami, FL.
- Hinton, G. E., & Anderson, J. A. (1981). *Parallel models of associative memory*. Hillsdale, NJ: Erlbaum.
- Hinton, G. E. (1981). Implementing semantic networks in parallel hardware. In G. E. Hinton & J. A. Anderson (Eds.), *Parallel Models of Associative Memory*. Hillsdale, NJ: Erlbaum.
- Hinton, G. E. & Sejnowski, T. J. (1983). Optimal Perceptual Inference. *Proceedings of the IEEE*

- conference on Computer Vision and Pattern Recognition.* Washington, DC: IEEE Computer Society.
- Hinton, G. E., & Sejnowski, T. J. (1986). Learning and relearning in Boltzmann machines. In D. E. Rumelhart, J. L. McClelland, & the PDP research group (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition.* Cambridge, MA: Bradford Books.
- Hinton, G. E., McClelland, J. L. & Rumelhart, D. E. (1986). Distributed representations. In D. E. Rumelhart, J. L. McClelland, & the PDP research group (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition..* Cambridge, MA: Bradford Books.
- Hinton G.E. (1986). Learning distributed representations of concepts. *Proc. Eighth Annual Conference of the Cognitive Science Society.* Amherst, Mass..
- Hinton, G. E. (1987). Learning translation invariant recognition in a massively parallel network. *Proc. Conf. Parallel Architectures and Languages Europe.* Eindhoven, The Netherlands.
- Hinton, G. E., McClelland, J. L. & Goodhill, G. J. (1987). *Recirculation: A biologically plausible approximation to back-propagation* Manuscript in preparation). Carnegie-Mellon University, Pittsburgh PA,
- Hogg, T. & Huberman, B. A. (1984). Understanding biological computation: Reliable learning and recognition. *Proceedings of the National Academy of Sciences*, 81, 6871-6875.
- Holland, J. H. (1975). *Adaptation in Natural and Artificial Systems.* University of Michigan Press.
- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences U.S.A.*, 79, 2554-2558.
- Hopfield, J. J., Feinstein, D. I. & Palmer, R. G. (1983). "Unlearning" has a stabilizing effect in collective memories. *Nature*, Vol. 304.
- Hopfield J.J., & Tank D.W. (1985). "Neural" computation of decisions in optimization problems. *Biological Cybernetics*, 52, 141-152.
- Hummel, R. A., & Zucker, S. W. (1983). On the foundations of relaxation labeling processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence, PAMI-5*, 267-287.
- Kerszberg, M. & Bergman, A. (1986). The evolution of data processing abilities in competing automata. *Proc. Computer Simulation in Brain Science.* Copenhagen, Denmark.
- S. Kirkpatrick, Gelatt, C. D., & Vecchi, M. P. (1983). Optimization by simulated annealing. *Science*, 220, 671-680.
- Kohonen, T. (1977). *Associative memory: A system-theoretical approach.* Berlin: Springer.

- Kohonen, T. (1982). Clustering, taxonomy, and topological maps of patterns. M. Lang (Ed.), *Proceedings of the Sixth International Conference on Pattern Recognition*. Silver Spring, MD, IEEE Computer Society Press.
- Le Cun, Y. (1985). A learning scheme for asymmetric threshold networks. *Proceedings of Cognitiva 85*. Paris, France.
- Linsker, R. (1986). From basic network principles to neural architecture: Emergence of spatial opponent cells. *Proceedings of the National Academy of Sciences USA*, 83, 7508-7512.
- Linsker, R. (1986). From basic network principles to neural architecture: Emergence of orientation-selective cells. *Proceedings of the National Academy of Sciences USA*, 83, 8390-8394.
- Linsker, R. (1986). From basic network principles to neural architecture: Emergence of orientation columns. *Proceedings of the National Academy of Sciences USA*, 83, 8779-8783.
- Linsker, R. (1987). Development of feature-analyzing cells and their columnar organization in a layered self-adaptive network. R. Cotterill (Ed.), *Computer Simulation in Brain Science*. Cambridge University Press.
- Marroquin, J. L. (1985). *Probabilistic Solution of Inverse Problems*. Doctoral dissertation, MIT,
- Minsky, M. L. (1954). *Theory of neural-analog reinforcement systems and its application to the brain-model problem* Ph.D. Dissertation. Princeton University,
- Minsky, M. & Papert, S. (1969). *Perceptrons*. Cambridge, Mass: MIT Press.
- Minsky, M. (1977). Plain Talk about Neurodevelopmental Epistemology. *5th International Joint Conferences on Artificial Intelligence*, Vol 2. .
- Parker, D. B. (April 1985). *Learning-logic* (Tech. Rep.). TR-47, Sloan School of Management, MIT, Cambridge, Mass.,
- Parker, D. B. (1987). *Second order back-propagation: An optimal adaptive algorithm for any adaptive network* Unpublished manuscript), .
- Pearl, J. (1986). Fusion, propagation, and structuring in belief networks. *Artificial Intelligence*, 29(3), 241-288.
- Pearlmutter, B. A. & Hinton, G. E. (1986). G-Maximization: An unsupervised learning procedure for discovering regularities. J. S. Denker (Ed.), *Neural Networks for Computing: American Institute of Physics Conference Proceedings 151*. .
- Plaut, D. C. and Hinton, G. E. (1987). Learning sets of filters using back-propagation. *Computer Speech and Language*, .

- Prager, R., Harrison, T. D. and Fallside, F. (1986). Boltzmann Machines for speech recognition. *Computer Speech and Language*, 1, 1-20.
- Rosenblatt, F. (1962). *Principles of neurodynamics*. New York: Spartan Books.
- Rumelhart, D. E., & McClelland, J. L. ((in press) 1985). On the acquisition of the past tense in english. In McClelland, J. L. and Rumelhart, D. E. (Eds.), *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 2: Applications*. Cambridge, MA: Bradford Books.
- Rumelhart, D. E. and Zipser, D. (1985). Competitive Learning. *Cognitive Science*, 9, 75-112.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by back-propagating errors. *Nature*, 323, 533-536.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by error propagation. In D. E. Rumelhart, J. L. McClelland, & the PDP research group (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition*. Cambridge, MA: Bradford Books.
- Saund, E. (1986). Abstraction and representation of continuous variables in connectionist networks. *Proceedings of the Fifth National Conference on Artificial Intelligence*. Los Altos, California, Morgan Kauffman.
- Sejnowski, T. J., & Rosenberg, C. R. (1987). Parallel networks that learn to pronounce English text. *Complex Systems*, 1, 145-168.
- Simon, H. A. (1969). *The Sciences of the Artificial*. MIT Press.
- Sutton, R. S. (Feb 1984). *Temporal credit assignment in reinforcement learning*. Doctoral dissertation, University of Massachusetts, Amherst, MA, (COINS Technical Report 84-02).
- Sutton, R. S. (1987). *Learning to predict by the method of temporal differences* (Tech. Rep.). GTE Laboratories technical report TR87-509.1,
- Terzopoulos, D. (1984). *Multiresolution computation of visible surface representations* Ph.D. Dissertation). Department of Electrical Engineering and Computer Science, MIT, Cambridge MA,
- Tesauro, G. (1987). *Scaling relationships in back-propagation learning: Dependence ofn training set size* Manuscript). Center for Complex Systems Research, University of Illinois at Urbana-Champaign,
- Thatachar, M. A. L. & Sastry, P. S. (1985). *Learning optimal discriminant functions through a cooperative game of automata* Technical Report EE/64/1985). Department of Electrical Engineering, Indian Institute of Science, Bangalore-560012, India,
- Valiant, L. G. (1984). A theory of the learnable. *Communications of the ACM*, 27:11, 1134-1142.

- Volper, D. J. & Hampson, S. E. (1986). Connectionist models of Boolean category representation. *Biological Cybernetics*, 54, 393-406.
- Von der Malsburg, C. (1973). Self-organization of orientation sensitive cells in striate cortex. *Kybernetik*, 14, 85-100.
- Weisbuch, G. & Fogelman-Soulie, F. (1985). Scaling laws for the attractors of Hopfield networks. *Journal of Physics Letters*, 46, 623-630.
- Widrow, B. & Hoff, M. E. (1960). Adaptive switching circuits. *IRE WESCON Conv. Record, Part 4*.
- Williams, R. J. (1986). *Reinforcement learning in connectionist networks: A mathematical analysis* (Technical Report). Institute for Cognitive Science, University of California San Diego, La Jolla, CA.
- Willshaw, D. J., Von der Malsburg, C. (1979). A marker induction mechanism for the establishment of ordered neural mapping: Its application to the retino-tectal connections. *Phil. Trans. Proc. R. Soc. Lond.*, B, 287, 203-243.
- Willshaw, D. (1981). Holography, associative memory, and inductive generalization. In G. E. Hinton & J. A. Anderson (Eds.), *Parallel models of associative memory*. Hillsdale, NJ: Erlbaum.