

Gain Variation in Recurrent Error Propagation Networks

Steven J. Nowlan

Computer Science Department, University Of Toronto,
10 Kings College Road, Toronto, Ontario M5S 1A4, Canada*

Abstract. A global gain term is introduced into recurrent analog networks. This gain term may be varied as a recurrent network settles, similar to the way temperature is varied when "annealing" a network of stochastic binary units. An error propagation algorithm is presented which simultaneously optimizes the weights and the gain schedule for a recurrent network. The performance of this algorithm is compared to the standard back propagation algorithm on a difficult constraint satisfaction problem. An order of magnitude improvement in the number of learning trials required is observed with the new algorithm. This improvement is obtained by allowing a much larger region of weight space to satisfy the problem. The simultaneous optimization of weights and gain schedule leads to a qualitatively different region of weight space than that reached by optimization of the weights alone.

1. Introduction

Neural networks have received much attention recently as plausible models for studying the computational properties of massively parallel systems. They have been applied to a variety of problems in signal processing, pattern recognition, speech processing, and more traditional AI problems. These models developed from earlier work in associative memories and models of cooperative computation in early vision processing [3,5,19].

Learning algorithms have been developed [1,18] that enable these networks to learn internal representations, allowing them to represent complex non-linear mappings. Two distinct types of networks have been studied quite extensively. The first of these uses analog units with a sigmoidal I/O function [8], and an error propagation algorithm for updating the weights to minimize an error function [16,17]. Most of these studies have focused on strictly feed-forward networks. The second type of network employs stochastic binary

*The author is currently visiting the University of Toronto, while completing a Ph.D. at Carnegie-Mellon University. This research was supported by an NSERC Scholarship.

units and symmetric connections. From an initial state these networks approach a low temperature equilibrium, in which low energy states have high probability. The weights in such networks may be updated by examining the difference in statistics between the equilibrium with the input and output states clamped and the equilibrium with these states unclamped [1].

Recent work has shown some interesting relationships between these two distinct models. Peterson and Anderson [13] have developed a continuous approximation to the Boltzmann machine algorithm, in which the stochastic binary units of the Boltzmann machine are replaced with analog units whose states are mean field approximations to the average states of corresponding stochastic binary units at equilibrium. They have shown significant speedup in convergence and improved generalization for interesting problems [14]. Hopfield [9] has shown that for a certain class of statistical estimation problems, the statistical network and the analog network have very closely related properties and learning algorithms. Provided that four conditions are met, the error propagation update rule for the weights in an analog feedforward network is a mean field approximation to the update rule for a statistical network using the Boltzmann machine algorithm. These four conditions are: the analog network must use an asymmetric divergence [6] rather than the more common mean square error function; the statistical averaging in the two state network is performed over the hidden and output units, but not the input units; the two networks have a small number of outputs; and the networks have only a single layer of non-interconnected hidden units.¹

This work explores another parallel between statistical and analog networks. Recurrent analog networks often show better convergence if a global gain term is introduced which may be varied over a single settling [7]. The result is a procedure similar to simulated annealing [11]. An error propagation scheme is presented which allows an analog network to "learn" its own gain schedule, and experimental results for a constraint satisfaction task show an order of magnitude speedup in learning when this approach is used.

2. Error propagation in recurrent analog networks

The recurrent analog networks that are considered here use a synchronous updating procedure, and place no restrictions on the nature of the connectivity matrix. The activation level of a unit j at time t is a nonlinear function of the input to the unit at time t :

$$y_{j,t} = g(x_{j,t}) \quad (2.1)$$

One possibility for g , used in our simulations, is the logistic function, $g(x) = 1/(1 + e^{-x})$. The input is a weighted sum of the states of units in the previous time step:

$$x_{j,t} = G_t \sum_i w_{ji} y_{i,t-1} \quad (2.2)$$

¹This allows hidden unit interactions to be ignored when using a first-order approximation.

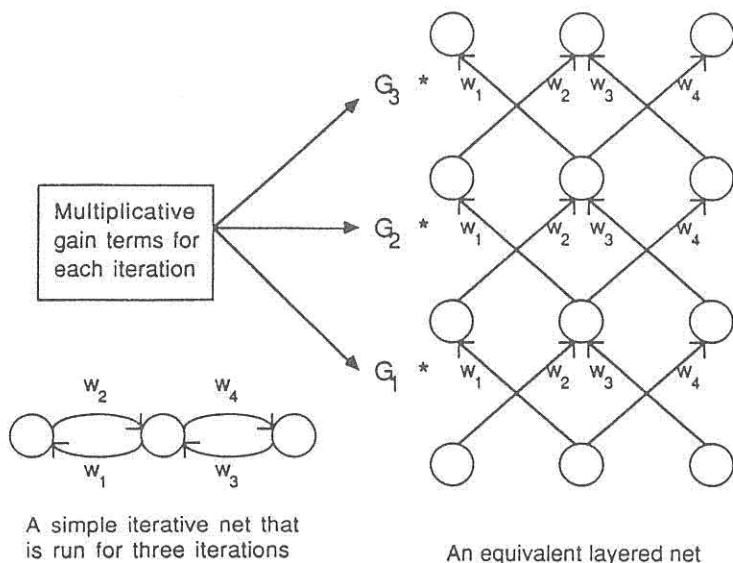


Figure 1: A recurrent network and the equivalent layered network. Corresponding weights in layers must be identical (i.e. w_1 has the same value in all layers), but the gain terms (G_i) vary between layers. Each layer corresponds to the state of the recurrent network at a different point in time.

The term G_t is a global gain term which premultiplies the input for every unit. There are also two distinguished subsets of units; \mathcal{I} the set of input units and \mathcal{O} the set of output units. The states of units in \mathcal{I} are determined by the environment.

The state vector for the recurrent network at time t can be treated as the vector of states of the t^{th} layer in an equivalent layered network (figure 1).

The trajectory of the state vector for the recurrent network is then represented by the state vectors of successive layers in the equivalent network. Since the weights in the recurrent network are fixed during the settling of the network, the set of weights between successive layers in the layered representation must be identical.

The dynamics for these networks can be expressed by finding the continuous differential equations equivalent to the discrete difference equations (2.1) and (2.2). This produces the following set of coupled differential equations:

$$\frac{dy_j}{dt} = -y_j + g(x_j) + I_j \quad (2.3)$$

where

$$x_j = G \sum_i w_{ji} y_i \quad (2.4)$$

I_j is defined to be zero for all units not in \mathcal{I} , and the time constant τ is assumed to be normalized to one. Equation (2.3) is a simple transformation of

$$\frac{dx_j}{dt} = -x_j + \sum_i w_{ji} g(x_i) + K_j \quad (2.5)$$

This equation has been studied by both Amari [2] and Hopfield [7]. Amari showed that in randomly connected networks with the dynamics of (2.5) the attractors were either stable or bistable.

When the attractors are fixpoints, one can use such networks to perform constraint satisfaction searches. The activities of units encode the values of the parameters of a problem, and the weights on the connections encode the constraints between the parameters. This approach has been used for classic optimization problems [8] and for parsing [20]. Coding the constraints into the weights by hand becomes a formidable task for large problems. This raises the possibility of devising learning algorithms which will manipulate the weights of a recurrent network to model the constraints of a specific problem through some training procedure.

There are several ways in which to pose the problem of modifying the weights. One approach would be to consider the fixpoint of the network for a specific input, and compute some error measure based on the distance of this fixpoint from a desired fixpoint. One could then use gradient descent to modify the weights to minimize this error measure [15]. An alternative is to consider not just the fixpoint, but the entire trajectory of the network. This approach was first suggested by Rumelhart, Hinton, and Williams [18] and is the approach taken here.

Consider a trajectory of length k for a recurrent network. There is an equivalent layered representation of this trajectory which has k distinct layers (figure 1). The standard back-propagation algorithm may be applied to this layered representation, if we define an error measure E for the final states of units in \mathcal{O} . Using this approach we can derive partials for the weights:

$$\frac{\partial E}{\partial w_{ji}} = \sum_t \frac{\partial E}{\partial x_{j,t}} G_t y_{i,t-1} \quad (2.6)$$

Note that the back-propagation proceeds through the sequence of states in the trajectory, and in particular that the partial of E with respect to w_{ji} will vary along that trajectory. Since the weights are fixed over the trajectory, some form of time averaging is required. Equation (2.6) uses a uniform time averaging, although versions which favor the terms near the end of the trajectory could also be used.² Thus the weights are being updated based on the average derivative over the trajectory. The disadvantage of using this

²The G_t terms actually do weight the derivatives, so the time averaging is not truly uniform.

approach to modifying the weights in the network is that it becomes necessary to store the entire trajectory for the back propagation phase. However, the advantage is that the network can be trained not just to have certain limit behavior, but also to have certain behavior along the trajectory followed to the limit. One obvious example is to force the network to learn fixpoints as attractors, by penalizing bistable behavior during the last few states of the trajectory. To allow control over the trajectory, the back propagation procedure is modified slightly to allow directly observed error terms in time steps other than the last to be added to the back propagated terms for units in \mathcal{O} . (This is equivalent to specifying desired states for intermediate layers in a layered network.)

3. Variable gain in recurrent analog networks

The term G_i in equation (2.6) determines the steepness of the nonlinearity in the recurrent network, and has been referred to as the system gain [8]. This term may be allowed to vary during the settling of the network. For example, an increasing gain in a network with mutually inhibitory connections can implement a winner-take-all network that converges quickly. Hopfield and Tank [8] found that increasing the gain slowly as their analog network settled increased the quality of the solution found by the network. They suggested that increasing the gain in this fashion was analogous to following the effective field solution from a high temperature, resulting in a final state near the thermodynamic ground state of the system.

An energy function, a non-increasing function of the state vector along any trajectory of the system, is a very useful tool for understanding the dynamic behavior of a recurrent analog network. For symmetric networks with the dynamics defined by (2.3), Hopfield [7] suggested an energy function of the form

$$\mathcal{E} = -\frac{1}{2} \sum_j \sum_i w_{ji} y_i y_j - \sum_i K_i y_i + \sum_i \int g^{-1}(y_i) dy_i \quad (3.1)$$

where $K_i = g^{-1}(I_i)$. To show \mathcal{E} is an energy function we note that $\frac{d\mathcal{E}}{dt} = \sum_i \frac{\partial \mathcal{E}}{\partial y_i} \frac{dy_i}{dt}$, so it is sufficient to show $\frac{\partial \mathcal{E}}{\partial y_i} \frac{dy_i}{dt} \leq 0$ for any i . Taking the partial of (3.1) with respect to y_i and combining with (2.3) yields

$$\frac{\partial \mathcal{E}}{\partial y_i} \frac{dy_i}{dt} = - \left(\left[\sum_j w_{ji} y_j + g^{-1}(I_i) \right] - g^{-1}(y_i) \right) \left(\left[g \left(\sum_j w_{ji} y_j \right) + I_i \right] - y_i \right) \quad (3.2)$$

The first parenthesized term in (3.2) is the difference between the current input to unit i and the input required to produce the current state of i . The second term is the difference between the state produced by the current input to unit i and the actual state of unit i . Provided that the function g

is monotonically increasing, the two parenthesized terms in (3.2) must be of the same sign, so $\frac{\partial \mathcal{E}}{\partial y_i} \frac{dy_i}{dt} \leq 0$ for any i . In addition to being monotonically increasing, g must also be invertible, and the inverse must be integrable.

The logistic function meets these conditions, and in this case both g^{-1} and its integral have simple forms:

$$g^{-1}(y) = \ln y - \ln(1 - y)$$

$$\int g^{-1}(y) dy = y \ln y - (1 - y) \ln(1 - y)$$

In the high gain limit when y_i is 0 or 1, g^{-1} goes to zero and (3.1) is identical to the energy function for stochastic networks proposed by Hopfield [10].

Although (3.1) has only been shown to be an energy function for symmetric networks, we have found that it is also a strictly decreasing function of the trajectory for all of the asymmetric networks we have simulated. This may be because these asymmetric networks tend to be very nearly symmetric once the weights have been learned.

Kirkpatrick [11] suggested that “simulated annealing” in stochastic networks is effective because of two features. First, it allows a system to occasionally take uphill steps, to allow it to escape from local minima. Second, it has the effect of averaging many solutions at high temperatures, which tends to smooth the energy surface. If the energy surface is reasonably well behaved, this smoothing can produce a multiresolution search. At high temperatures regions with many energy minima will appear as smooth valleys, while local minima in regions of high average energy will be smoothed out. Initially the system will head towards regions of low average energy, and as the temperature is lowered will settle into a local minima that is also close to a global minima.

The state evolution of a recurrent analog network can be regarded as a mean field approximation of the evolution of a stochastic network [8]. In this approximation the gain takes the role of an inverse temperature, and states of the analog network at low gain correspond to averages of states of the stochastic network at high temperature. This suggests that the energy surface of the analog network should be smoother with low gain and exhibit more local minima as the gain is increased. This effect may be seen clearly in figure 4 which shows the evolution of the energy surface for the 5 queens problem³ as the system gain is increased from 0.5 to 1.1. A multiresolution search can be carried out by the analog network by slowly increasing the gain as the network settles. This should result in a fixpoint of relatively low global energy.

Rather than determining a gain schedule in advance, an error propagation scheme can be used to decide how G_t should vary over the trajectory of the network. G_t is optimized by performing gradient descent in the error measure:

³This problem is discussed in the next section, and this energy surface is discussed in more detail.

$$\frac{\partial E}{\partial G_t} = \sum_j \sum_i \frac{\partial E}{\partial x_{j,t}} w_{ji} y_{i,t-1} \quad (3.3)$$

Once again averaging is necessary, in this case over all units in the network, since G_t is a global premultiplier.

Equations (2.6) and (3.3) can be combined to produce an algorithm which trains a network to exhibit desired trajectories by simultaneously modifying the weights and gain schedule. The trajectories learned correspond to a constraint satisfaction search via relaxation.

4. Performance on a constraint satisfaction problem

The author investigated the performance of the gain variation error propagation algorithm through simulations on some small problems [16]. For simple coding and sequencing tasks a network that learns a variable gain schedule learns to solve a problem several times faster than a similar network with fixed gain. In addition, the learned gain variation schedules outperformed several hand designed schedules on the same tasks. These tasks are all expressed in terms of I/O mappings; a prespecified output was required for each input. This makes the dynamics to be learned quite a bit easier. It is possible for the network to learn a trajectory from each input to the desired output without constructing a true attractor for that output.⁴ A task in which it is necessary to construct robust attractors which represent the problem constraints provides a much richer domain in which to study the performance of gain variation.

The problem selected is the n queens problem. This is a classical constraint satisfaction problem that was studied extensively by early AI researchers [4,12]. The general problem is to place n queens on an n by n grid of squares, such that there is no vertical, horizontal, or diagonal line through the grid that contains more than 1 queen. This problem is easily mapped into a network: Each cell of the grid is represented by a unit in the network, and each unit is fully connected to every other unit including itself.⁵ In addition, each unit has an external input line to carry environmental input. In this special case every unit is both an input and an output unit ($\mathcal{I} = \mathcal{O}$). One nice feature of this problem is that it can be easily scaled for larger values of n .

Given a random initial state vector, the network is required to settle into a final state which represents a valid solution to the n queens problem. A valid solution is one in which n units are above the *on level* (0.9), and all the rest are below the *off level* (0.1). In addition, no two *on* units can lie on the same vertical, horizontal, or diagonal line. To solve this problem, the network must construct stable attractors for the valid solutions to the problem and the set of attraction basins must span the entire input space.

⁴Additional simulations showed that this was in fact the case for several of the experiments discussed.

⁵The network must learn which of these connections are really needed to solve the task.

Since this task differs from the typical I/O mapping tasks given to error propagation algorithms, some care must be taken in deciding on an error measure. Given a random initial input state, there are in general many final states which are equally acceptable as solutions. One possibility is to measure the distance of the actual output from each of these final solutions, and take the minimum distance as the error to be propagated. This results in a form of nearest neighbor error measure.

An even more sophisticated training method, a form of shaping, may be used to produce good results and reduce the training time. The network is first presented with noisy versions of solution vectors as input. A solution is chosen at random, and then noise uniformly distributed between 0 and η (initially 0.2) is added to units that are off and subtracted from units that are on. This noisy vector is clamped to the inputs, and the network allowed to settle for γ cycles. During the last three cycles, the mean square distance between the output vector and the solution vector used to generate the input is calculated as the error measure. The error is taken over the last three cycles to force the network to learn attractors which are fixpoints. The network is trained in this fashion until its average error (normalized by the vector lengths) is less than 20 percent over the last 50 trials. At this point, the same training regime is continued, except that the input is presented during the first cycle only, rather than being clamped. During this second phase η is gradually increased to 0.4. Once the average error over 50 trials is less than 10 percent, a final training phase is performed in which initial states are randomly generated, and the nearest neighbor error measure is used over the last three trials. The first phase of training establishes the attractors, the second phase stabilizes the attractors independent of external input, and the third phase ensures robustness of the attractors.

The simulation results are summarized in table 1. One obvious anomaly in the table is the difference between the 5 queens and 6 queens problems. The gain variation technique shows a clear advantage for the 5 queens problem (an order of magnitude improvement), but the performance with and without gain variation is nearly identical for the 6 queens problem. The answer lies in the second column of the table. There are only 4 solutions for the 6 queens problem, and two of these solutions are simple mirror images of the other two. It is quite easy for the network to set the unit biases to favor the small set of units which appear in these 4 solutions. This "trick" makes establishing stable attractors quite easy. On the other hand for the 5 queens problem, and the larger problems, there are sufficient solutions so that almost every unit is active in at least one solution, so the biases alone cannot be used to give the network a head start. Under these circumstances, the dynamical behavior of the network becomes much more important, and so a much stronger advantage is shown by the gain variation algorithm.

The set of weights learned by the network for one example of the 5 queens problem (figure 2) shows that in its solution the network has extracted the essential constraints of the task. Each unit has learned to develop a positive weight to itself, and negative weights along all horizontal, vertical and

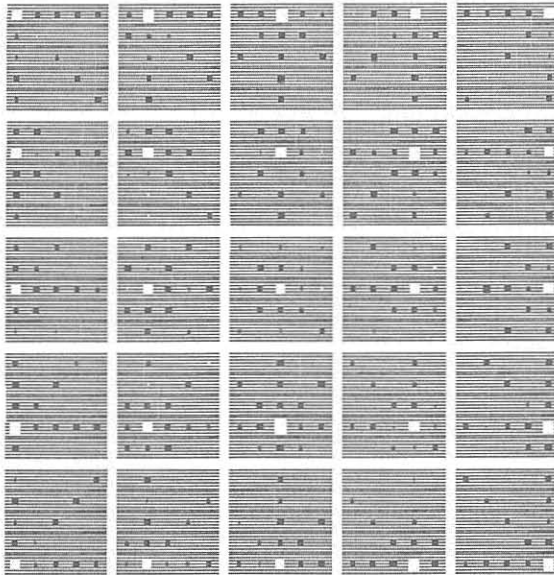


Figure 2: Weights learned for the 5 queens problem. The weight display is recursive, each large grey square represents one unit. Within each large grey square is a 5 by 5 grid of squares which represent the weight of the connection from that unit to every other unit in the network. Black squares represent negative weights, white squares represent positive weights, and the size of the square represents the magnitude of the weight. Weight decay was applied to drive all non-essential weights to zero, to simplify the weight display.

n	Solutions	No Gain Variation				Gain Variation			
		Phase 1	Phase 2	Total	ϵ	Phase 1	Phase 2	Total	$\gamma\text{-}\epsilon$
4	2	115	420	677	0.1	70	232	350	1.0×10^{-3}
5	10	6371	11048	19820	0.01	371	1048	1540	1.0×10^{-4}
6	4	371	1048	1420	0.1	397	698	1242	1.0×10^{-3}
7	40	20600	57480	83200	0.01	1280	4160	6080	1.0×10^{-4}
8	92	50000	—	50000	0.01	1960	12600	16200	1.0×10^{-4}

Table 1: Simulation results for various sizes of the n queens problem. ϵ is the size of the weight step; $\gamma\text{-}\epsilon$ is the size of the gain update. The numbers under the columns Phase 1, Phase 2, and Total are the number of weight updates performed in each training phase (see text) and in total. One update was performed after every 20 training examples. The results for the 4 queens problem are averaged over 50 runs, for the 5 and 6 queens problems the averages are over 20 runs. Only 2 runs are reported for the 7 queens, and 1 run for the 8 queens. The algorithm with no gain variation was not able to meet the 20 percent error criteria for phase 1 for the 8 queens problem and was terminated after 50000 updates.

diagonal lines which the unit lies on. This means that all of the units have a natural tendency to turn on, and along each line (in any orientation) a winner-take-all network has formed, so that the stable states are those in which only one unit on a particular line is on. This is a very natural representation of the original problem constraints, which stipulated that no two queens could lie on the same horizontal, vertical or diagonal line.

The dynamic behavior of a 5 queens network can be seen in figure 3. In the third of the five examples we can see the network placing a queen in the third row when the initial configuration contains no dominant unit in this row. In the fifth example, we can see competition between two initially dominant units in both the first and second rows, and also the creation of a dominant unit in the third row. The performance of the network is quite robust, even on highly ambiguous inputs.

The energy surface for one simulation of the 5 queens problem (figure 4) gives further insight into the nature of the solution learned. A cross section of the energy surface for the states of two units (all other units are fixed at a solution point) is shown for several different values of system gain. The two units shown are adjacent units on the first row of the network. At this particular solution point there are no units yet “on” in the first row, so it is desired that one of these two units should turn on, but not both of them. In terms of the problem constraints, the most desirable fixpoints are the corners (0.0,1.0) or (1.0,0.0). The origin would also be acceptable, but the corner (1.0,1.0) would violate a constraint. With this background, we can now interpret figure 4. Consider first part (d) of the figure, the high gain case. There are 4 attractors on this surface, represented by the 4 major

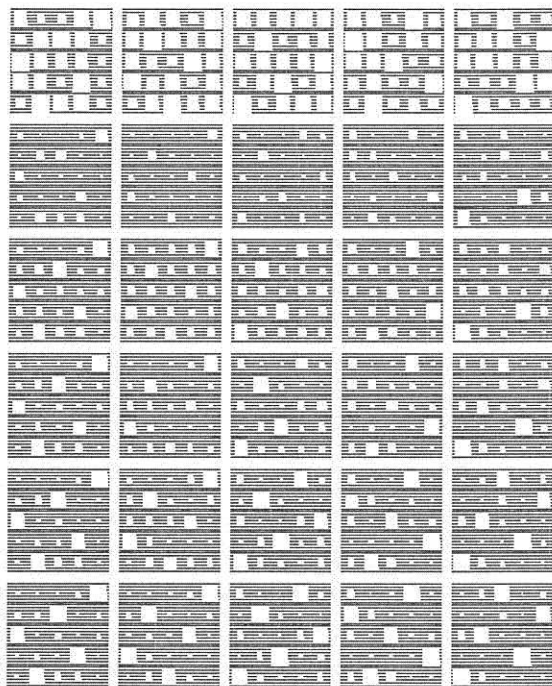
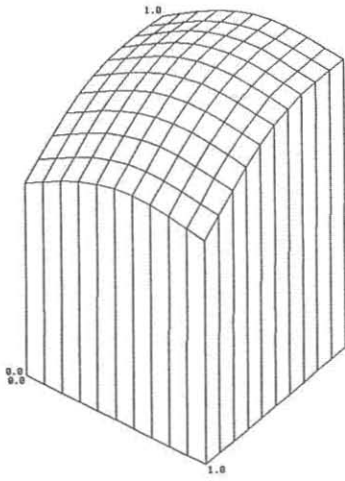
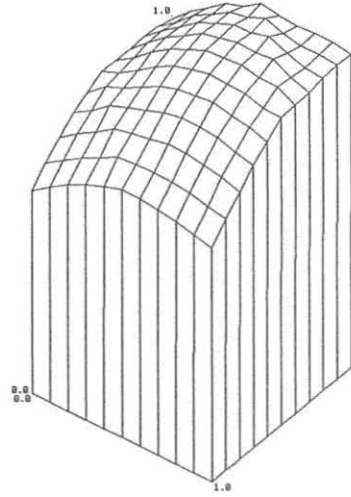


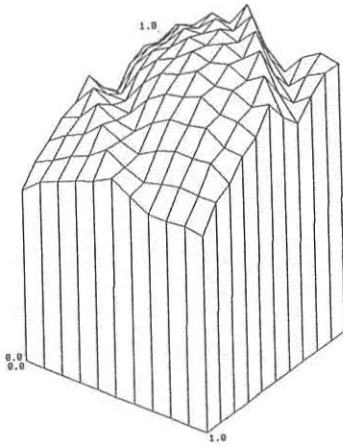
Figure 3: Display of the dynamics of the 5 queens problem for 5 different initial states. Each column represents the settling of the network for one case. The large black squares contain the activity levels of all 25 units in the network, arranged in a 5 by 5 grid. The size of the white square in each grid position is proportional to the activity of the corresponding unit. In all cases the state reached at the end of 6 cycles is a stable fixpoint.



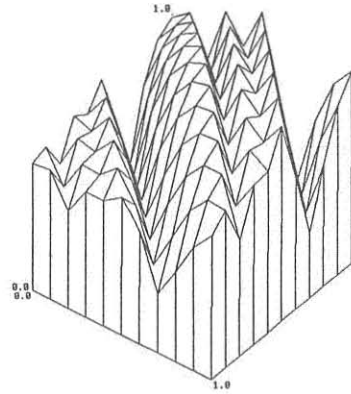
(a) 5 Queens Energy Surface:
Gain 0.5



(b) Queens Energy Surface:
Gain 0.7



(c) 5 Queens Energy Surface:
Gain 0.9



(d) 5 Queens Energy Surface:
Gain 1.1

Figure 4: A cross section of the energy surface of the 5 queens problem for one set of weights and 4 different values of gain. The x and z axes are the states of two units in the first row of the network, which range from 0.0 to 1.0.

valleys. (There is also a partial valley between the third and fourth valleys, but this will spill into the third valley, and is not an attractor.) The second and third valleys correspond to attractors that are near the corners $(0.0, 1.0)$ or $(1.0, 0.0)$. The lowest points of these valleys are at the two ends, when one unit is fully off and the other is mostly on. The first valley represents an attractor near the origin; since this solution is less preferred, this attractor is not as deep. The most interesting feature of this energy surface is the fourth valley. This is a strong attractor near the corner $(1.0, 1.0)$. This is a spurious attractor, which violates the constraints for the first row. The system is able to avoid being trapped by this attractor because it does not settle with constant gain. We can see in figure 4 (a) that at low gain the energy surface slopes smoothly away from the corner $(1.0, 1.0)$. A network started near this corner at low gain will move away from this corner, and will be beyond the collection basin for this attractor as it begins to form with increasing gain (figure 4 b and c). The shape of the error surface, with its spurious attractor, suggests that the set of weights found in this simulation would not be a solution if the network used a constant gain schedule. The experiments mentioned in the discussion at the end of this paper support this observation.

Our experiments with the simple I/O mapping networks showed a tendency for the learned gain schedules to be "annealing" schedules, starting at a low gain and increasing it as the network settled [16]. This same effect is observed in the simulations for the n queens problem (figure 5). Here there is also another effect. The gain increases steadily until the network reaches its fixpoint (after five time steps in the figure) and then declines slightly.

5. Discussion

Varying the gain for a recurrent network as it settles has been suggested elsewhere [8,15], as has an analogy between gain in recurrent analog networks, and temperature in statistical networks [2,15]. The unique aspect of this work is the use of an error propagation scheme to *simultaneously* optimize the weights and gain schedule for a recurrent network. The empirical results presented here show that for constraint satisfaction problems in which a complex attractor structure must be developed, the parallel optimization of the weights and gain schedule can produce an order of magnitude speed-up in the convergence of the optimization. What is not clear is whether this speed-up is obtained by following a shorter path to the same weight region that would be reached by optimizing the weights alone, or whether a qualitatively different region of weight space is reached by the combined optimization.

Some simple simulations tend to support the latter hypothesis. If an n queens network that has been trained with gain variation has its gain schedule modified so that the gain is a constant value (for example, the mean of the gain schedule), qualitatively different behavior is observed from the network. Additional stable states are observed, which do not correspond to solutions to the n queens problem, and which are not stable when the net-

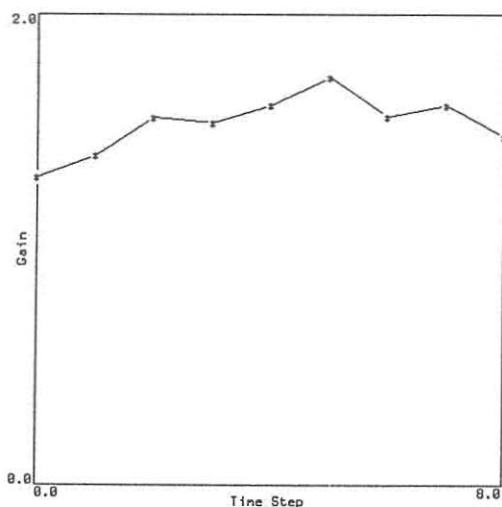


Figure 5: Graph of the gain schedule learned for one run of the 5 queens problem. Qualitatively the graphs for all other runs were identical, exhibiting an initial section of steady increase, and a slight drop of gain after 5 or 6 time steps.

work is allowed to settle using the gain schedule. This behavior is expected. Annealing the analog network with the gain schedule will force the network into states nearer the thermodynamic ground state, and away from higher energy fixpoints. It would appear that the use of gain variation allows the network to find a set of weights that have an attractor structure with many high energy spurious attractors, in addition to the low energy attractors that correspond to solutions to the task. It is apparently much easier to find a set of weights with an attractor structure of this form, rather than attractors which correspond only to task solution points. This leads to the hypothesis that the speed-up in convergence is obtained by allowing a much larger region of weight space to satisfy the problem, and that the combined optimization leads to a qualitatively different region of weight space than by optimization of the weights alone.

There is an additional factor which may account for some of the speed-up in learning observed with the parallel optimization of the weights and gain schedule.⁶ The error surface for many problems that back propagation is applied to is characterized by ravines with steep sides in most directions, but a shallow descent in one direction. Once the weight vector is aligned with the floor of the ravine, one can move quite rapidly along this floor by simply adjusting a global gain term. To examine this effect, several of the 5 queens simulations were repeated with a gain term that was learned by

⁶This was suggested by Geoff Hinton and Mike Mozer, personal communication.

error propagation, but was constrained to be constant during a settling (as the weights were constrained). The average total number of weight updates (see table 1) over 15 runs was 8371, compared to 1540 for the variable gain algorithm, and 19820 for the algorithm with no gain. Similar trends appeared in the 6 and 7 queens problems. The fixed gain algorithm is a factor of two to three faster than the simple back-propagation algorithm, but still five to six times slower than the parallel optimization of the weights and gain schedule, suggesting that the ability to scale all of the weights accounted for a small part of the speed-up observed in the n queens problem. Nevertheless, the speed-up was significant enough to suggest using a similar scale factor in layered networks.

References

- [1] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski, "A learning algorithm for Boltzmann machines," *Cognitive Science*, **9** (1985) 147-169.
- [2] Shun-Ichi Amari, "Characteristics of random nets of analog neurons," *IEEE Transactions on Systems, Man, and Cybernetics*, **2**(5) (1972) 643-657.
- [3] J. A. Anderson and G. E. Hinton, "Models of information processing in the brain" in *Parallel models of associative memory*, J. A. Anderson and G. E. Hinton, eds. (Erlbaum, Hillsdale, NJ, 1981).
- [4] Edward A. Feigenbaum and Avron Barr, eds., *The Handbook of Artificial Intelligence*, **1** (Pitman, London, 1981).
- [5] J. A. Feldman and D. H. Ballard, "Connectionist models and their properties," *Cognitive Science*, **6** (1982) 205-254.
- [6] Geoffrey E. Hinton, *Connectionist Learning Procedures*, Department of Computer Science CMU-CS-87-115, Carnegie-Mellon University, Pittsburgh, PA, June 1987.
- [7] J. J. Hopfield, "Neurons with graded response have collective computational properties like those of two-state neurons," *Proceedings of the National Academy of Sciences U.S.A., Bio.*, **81** (1984) 3088-3092.
- [8] J. J. Hopfield and D. W. Tank, "'Neural' computation of decisions in optimization problems," *Biological Cybernetics*, **52** (1985) 141-152.
- [9] J. J. Hopfield, "Learning algorithms and probability distributions in feed-forward and feed-back networks," *PNAS*, **84** (1987) 8429-8433.
- [10] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," *Proceedings of the National Academy of Sciences U.S.A.*, **79** (1982) 2554-2558.
- [11] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *Science*, **220** (1983) 671-680.

- [12] Nils J. Nilsson, *Principles of Artificial Intelligence*, (Tioga, Palo Alto, CA, 1980).
- [13] C. Peterson and J. R. Anderson, *A Mean Field Theory Learning Algorithm for Neural Networks*, MCC Technical Report EI-259-87, Microelectronics and Computer Technology Corporation, Austin, TX, August 1987.
- [14] C. Peterson and J. R. Anderson, *Neural Networks and NP-complete Optimization Problems: A Performance Study on the Graph Bisection Problem*, Technical Report MCC-EI-287-87, Microelectronics and Computer Technology Corporation, Austin, TX December 1987.
- [15] F. J. Pineda, "Generalization of back propagation to recurrent and higher order neural networks" in *Proceedings of IEEE Conference on Neural Information Processing Systems*, IEEE, Denver, CO, November 1987.
- [16] D. C. Plaut, S. J. Nowlan, and G. E. Hinton, *Experiments on Learning by Back Propagation*, Department of Computer Science CMU-CS-86-126, Carnegie-Mellon University, Pittsburgh, PA 1986.
- [17] D. C. Plaut and G. E. Hinton, "Learning sets of filters using back-propagation," *Computer Speech and Language*, (1987).
- [18] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by back-propagating errors, *Nature*, **323** (1986) 533-536.
- [19] D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, eds. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Volume I (MIT Press, Cambridge, MA, 1986).
- [20] B. Selman and G. Hirst, "Parsing as an energy minimization problem" in *Genetic Algorithms and Simulated Annealing*, Lawrence Davis, ed. (Pitman, London, 1987) 155-168.