

On-line learning in the committee machine

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

1995 J. Phys. A: Math. Gen. 28 1615

(<http://iopscience.iop.org/0305-4470/28/6/016>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 171.66.16.68

The article was downloaded on 02/06/2010 at 01:57

Please note that [terms and conditions apply](#).

On-line learning in the committee machine

Mauro Copelli† and Nestor Caticha‡

Instituto de Física, Universidade de São Paulo, CP 20516, 01498 São Paulo, SP, Brazil

Received 24 October 1994, in final form 19 January 1995

Abstract. The dynamics of learning from examples in the $K = 3$ non-overlapping committee machine with single presentation of examples is studied. The optimal algorithm, in the sense of mean generalization, is obtained from a variational analysis of the differential equations which describe the dynamics. The agreement of the theoretical predictions and the results of numerical simulations is excellent. The optimized dynamics has the extra advantage with respect to the non-optimized cases in that it uncouples the differential equations which describe the evolution of the relevant parameters, i.e. the student–teacher overlap and the norm of the student synaptic vector. This, in turn, translates into the possibility of constructing useful practical optimized on-line algorithms. For the optimal algorithm the generalization error decays as $\sim 0.88\alpha^{-1}$, the same nominal error as for the simple perceptron with optimized dynamics.

1. Introduction

The application of the methods of statistical mechanics to the study of learning and generalization in perceptrons with no hidden units has yielded a wealth of results which lead, in a natural way, to the study of more realistic and complex nets. The feedforward nets with a layer of hidden units connected through fixed weights to the output layer have been the natural candidates to continue the study of neural networks.

The supervised learning process can typically be cast as a problem of minimizing some energy function constructed from a set of examples. This minimization leads naturally to a learning dynamics. If this minimization is performed in the presence of noise then the problem is that of statistical mechanics at finite temperature. Thus the learning algorithms describe a sequence of steps which tend to minimize a free energy. In the iterated form, these algorithms extract all the information contained in the *a priori* given energy function. From a practical point of view, this might be a computationally expensive operation. For the perceptron with no hidden units it has been shown that if the dynamics is properly chosen, then the first step of the minimization leads to an asymptotic decay of the generalization errors that is almost as good as the exhaustive iterated learning. For example, the generalization error in the case of the Boolean perceptron with real weights decays as $0.88\alpha^{-1}$, compared to the Bayesian bound of Oppen and Haussler [12] of $0.44\alpha^{-1}$. In the case of learning by queries, by judiciously choosing the energy function, the error can be made to decay exponentially fast [7].

The main point of the above discussion is that by correctly choosing the dynamics, the computational effort of learning can be drastically reduced. The on-line learning by single presentation of examples, which lately has been studied by several authors [10, 8, 5, 1], can

† E-mail address: copelli@if.usp.br

‡ E-mail address: nestor@if.usp.br

thus lead to very efficient learning algorithms. In the case of nets which are not very useful in real-life applications this might not be very important. This is not the case for nets with hidden units, where any reduction in the training times or improvement of performance will certainly be welcomed, and which are applicable in a large number of real problems.

The object of this paper is to study the single presentation of example algorithms which lead to the best possible generalization (in the mean) in the committee machine with a non-overlapping architecture. The method generalizes to nets with hidden units—see [8].

In section 2 the model is presented, an expression for the generalization error as a function of the overlaps is obtained and the equations governing the evolution of the order parameters for a general algorithm are deduced. In section 3 we obtain the optimal weight function through a variational argument. This leads to an algorithm which belongs to the class of the 'expected stability' algorithms [8]. In section 4 the theoretical results, obtained by numerical integration of the evolution equations are compared to the results of numerical simulations. Section 5 contains some concluding remarks.

2. The generalization error and the learning dynamics

Since we are interested in the problem of generalization or rule extraction within the framework of supervised learning, we build a learning set with the help of a teacher network [3, 4], which for simplicity we take to have the same architecture as the student net.

The K non-overlapping committee we deal with is a set of K independent Boolean perceptrons or branches with N/K inputs units each. The notation we use is such that every N -dimensional vector V can be thought of as K branch-vectors $V = (V_1, \dots, V_K)$. We will restrict ourselves to the $K = 3$ case.

The learning set is a set of P pairs $\{(S, \sigma_B)\}$ where $S = (S_1, S_2, S_3)$ with $S_{kj} = \pm 1$, $k = 1, 2, 3$, $j = 1, \dots, N/3$ and σ_B is the teacher output (see equation (1) below). The synaptic weights of the teacher are denoted by $B = (B_1, B_2, B_3)$ and the normalization $\sum_{j=1}^{N/3} B_{kj}^2 \equiv B_k^2 = 1$ can be imposed without loss of generality.

Upon presentation of an input vector each branch perceptron gives a partial output

$$\sigma_{Bk} = \text{sign}(b_k) \quad k = 1, 2, 3$$

where $b_k \equiv B_k \cdot S_k$. The set $\{\sigma_{Bk}\}$ is the unaccessible internal representation of S in the teacher net. The teacher's final output is made up from the internal representation

$$\sigma_B = \text{sign}(B) \quad (1)$$

where $B \equiv \sum_{k=1}^3 \sigma_{Bk}$. The student net is defined by a vector of real connections J and an output

$$\sigma_J = \text{sign}(H) \quad (2)$$

where $H \equiv \sum_{k=1}^3 \sigma_{Jk}$ and

$$\sigma_{Jk} = \text{sign}(h_k) \quad k = 1, 2, 3$$

$$h_k \equiv J_k \cdot S_k / |J_k|.$$

The aim of training the net is to obtain J , using information from the learning set, such that σ_J is near in some sense to σ_B . The meaning of 'near' here is in the sense of generalization, that is, we should find J which maximizes the probability of $\sigma_B = \sigma_J$ upon the presentation of a random input vector, uncorrelated with the learning set. That probability is called the generalization ability (g) and its complement is called the generalization error (e_g).

We deduce now the generalization error as a function of the overlaps $\rho_k = \mathbf{J}_k \cdot \mathbf{B}_k / J_k B_k$, where $J_k \equiv |\mathbf{J}_k|$. For given \mathbf{B} and \mathbf{J} , the generalization error is the average, over random examples, of

$$e_g^\mu = \frac{1 - \text{sign}(\mathcal{H}^\mu)\text{sign}(\mathcal{B}^\mu)}{2}$$

which is the error of classification of a single example (labelled by the μ index) (S_k^μ, σ_B^μ). The average over the examples can be written as

$$e_g = \int_{-\infty}^{+\infty} d\mathcal{H}^\mu \int_{-\infty}^{+\infty} d\mathcal{B}^\mu P(\mathcal{H}^\mu, \mathcal{B}^\mu) e_g^\mu = \int \int_{\mathcal{H}^\mu \mathcal{B}^\mu < 0} P(\mathcal{H}^\mu, \mathcal{B}^\mu) d\mathcal{H}^\mu d\mathcal{B}^\mu$$

where, supressing the μ index for simplicity,

$$P(\mathcal{H}, \mathcal{B}) = \int_{-\infty}^{+\infty} \frac{dx}{2\pi} e^{ix\mathcal{H}} \int_{-\infty}^{+\infty} \frac{dy}{2\pi} e^{iy\mathcal{B}} \prod_{k=1}^K (e_k \cos(x-y) + g_k \cos(x+y)) \quad (3)$$

and

$$g_k = 1 - e_k = 1 - \frac{\cos^{-1} \rho_k}{\pi} \quad (4)$$

which comes from the Gaussian correlated distribution of the internal fields of both the student and teacher branch perceptrons, as shown by Mato and Parga [11].

Defining $C_\pm = \cos(x \pm y)$, we have

$$P(\mathcal{H}, \mathcal{B}) = \int_{-\infty}^{+\infty} \frac{dx}{2\pi} e^{ix\mathcal{H}} \int_{-\infty}^{+\infty} \frac{dy}{2\pi} e^{iy\mathcal{B}} [(e_1 e_2 e_3) C_-^3 + (e_1 e_2 g_3 + e_1 g_2 e_3 + g_1 e_2 e_3) C_-^2 C_+ + (e_1 g_2 g_3 + g_1 e_2 g_3 + g_1 g_2 e_3) C_+^2 C_- + (g_1 g_2 g_3) C_+^3].$$

Collecting terms with $\mathcal{H}\mathcal{B} < 0$ we get

$$e_g = (e_1 e_2 e_3) + \frac{1}{2}(e_1 + e_2 + e_3 - e_1 e_2 - e_1 e_3 - e_2 e_3). \quad (5)$$

This result expresses the geometrical relation between the generalization error and the overlaps.

The evolution of the overlaps during the learning process is obtained from the learning algorithm prescription of how the synaptic weights are modified upon presentation of a new example, which we take as the most general linear and synchronous dynamics:

$$J_{ki}(\mu + 1) = J_{ki}(\mu) \left(1 - \frac{\Omega_k(\mu)}{N} \right) + \frac{F_k(\mu) S_{ki}(\mu)}{N} \quad (6)$$

where $F_k(\mu)$ defines the algorithm and $\Omega_k(\mu)$ just rescales the norm of \mathbf{J}_k . In the $N \rightarrow \infty$ limit we have the following equations:

$$\begin{aligned} \frac{\Delta \rho_k(\mu)}{\Delta \alpha} &= \frac{\rho_k(\mu) F_k(\mu)}{J_k(\mu)} \left(\frac{b_k(\mu)}{\rho_k(\mu)} - h_k(\mu) - \frac{F_k(\mu)}{2K J_k(\mu)} \right) \\ \frac{\Delta J_k(\mu)}{\Delta \alpha} &= J_k(\mu) \left(\frac{F_k^2(\mu)}{2K J_k^2(\mu)} + \frac{F_k(\mu) h_k(\mu)}{J_k(\mu)} - \Omega_k(\mu) \right) \end{aligned}$$

where $\alpha = \mu/N$ is the number of presented patterns per adjustable network weight. Averaging over random examples and taking the mean values ρ_k and J_k as approximations (which will turn out to be very good) for $\rho_k(\mu)$ and $J_k(\mu)$ we get the mean evolution equations:

$$\frac{d\rho_k}{d\alpha} = \frac{\rho_k}{J_k} \int D\mu F_k(\mu) \left(\frac{b_k(\mu)}{\rho_k} - h_k(\mu) - \frac{F_k(\mu)}{2K J_k} \right) \quad (7)$$

$$\frac{dJ_k}{d\alpha} = J_k \int D\mu \left(\frac{F_k^2(\mu)}{2K J_k^2} + \frac{F_k(\mu)h_k(\mu)}{J_k} - \Omega_k(\mu) \right) \quad (8)$$

where $D\mu$ stands for the measure over random patterns.

The dynamics is described by $2K = 6$ variables $\rho_1, \rho_2, \rho_3, J_1, J_2, J_3$ and by the initial conditions. The optimization procedure of the next section will uncouple the ρ set from the J set, leading not only to better performance but also to a simpler to analyse dynamics.

3. Dynamics optimization

From equations (5) and (4) it can be seen that minimization of e_g is obtained by maximizing the increase in the overlaps ρ_k 's for each newly presented example. A simple variational method applied to (7) leads to the optimal weight F_k^* , given by

$$F_k^* = K J_k \left(\frac{b_k}{\rho_k} - h_k \right) \quad (9)$$

where the μ index will be dropped for simplicity. However, the values of $\{b_j\}$ are not known. Define a new weight \tilde{F}_k in which the inaccessible variable b_k is replaced by some unknown function f_k which may depend on anything but the b_j 's, i.e.

$$\tilde{F}_k = K J_k \left(\frac{f_k}{\rho_k} - h_k \right).$$

Inserting \tilde{F}_k into (7) we get

$$\frac{d\rho_k}{d\alpha} = \rho_k \int D\mu \left(\frac{f_k b_k}{\rho_k^2} - \frac{f_k^2}{2\rho_k^2} - \frac{h_k b_k}{\rho_k} + \frac{h_k^2}{2} \right) \quad (10)$$

where we have used the shorthand

$$\int D\mu (\dots) = \sum_{\sigma_B = \pm 1} \int P(b_k | \sigma_B, \{h_j\}) P(\sigma_B, \{h_j\}) db_k \left[\prod_{j=1}^3 dh_j \right] (\dots).$$

After integrating on b_k , variational optimization on the f_k function leads to

$$f_k^* = \langle b_k \rangle_k$$

where the brackets $\langle \dots \rangle_k$ stand for the average of b_k given $\sigma_B, \{h_j\}$, i.e. average over the posterior distribution probability

$$f_k^* = \int_{-\infty}^{+\infty} db_k P(b_k | \sigma_B, \{h_j\}) b_k = \frac{\int_{-\infty}^{+\infty} db_k P(b_k, \sigma_B, \{h_j\}) b_k}{\int_{-\infty}^{+\infty} db_k P(b_k, \sigma_B, \{h_j\})}. \quad (11)$$

We have thus obtained the accessible optimal weight

$$\tilde{F}_k = K J_k \left(\frac{\langle b_k \rangle_k}{\rho_k} - h_k \right). \quad (12)$$

The meaning of this result is clear. The optimal weight that should be attached to a new example depends on b_k (equation (9)). This, however, does not comply with the rules of the game and the next best thing that can be done is to replace b_k by the expected value $\langle b_k \rangle_k$. Another point that might be raised is that the overlaps ρ_k 's are not accessible either. There are several ways out of this problem. The first and most obvious one is to obtain the value of $\rho_k(\alpha)$ by integration of the evolution differential equations. This leads to mean-field-type lower bounds on the generalization error. A second way out of the ρ -problem has been suggested for single-layer perceptrons [9], where a self-adaptive algorithm estimates

the value of ρ from a measure of the performance on the last few examples. We will now see that the optimization procedure will grant the possibility of a third solution, much more elegant and appealing than the previous two, especially from the point of view of applications.

In the special case where the examples are independent identical uniformly distributed random variables, and for the non-overlapping architecture which allows for branch factorization, we have

$$P(\{b_k, h_k\}) = \prod_{k=1}^3 P_0(b_k, h_k)$$

where

$$P_0(b_k, h_k) \equiv P_0(h_k)P_0(b_k|h_k) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}h_k^2\right) \frac{1}{\sqrt{2\pi(1-\rho_k^2)}} \exp\left(-\frac{(b_k - \rho_k h_k)^2}{2(1-\rho_k^2)}\right). \quad (13)$$

Taking advantage of the symmetry of the problem, we now calculate $P(\sigma_B, b_1, \{h_j\})$:

$$P(\sigma_B, b_1, \{h_j\}) = P_0(b_1, h_1) \prod_{j \neq 1} P_0(h_j) \int \delta\left(\sigma_B - \text{sign}\left(\sum_{i=1}^3 \text{sign}(b_i)\right)\right) \prod_{j \neq 1} P_0(b_j|h_j) db_j$$

where δ is the Kronecker delta-function. Introducing $\lambda_k \equiv \rho_k^{-1} \sqrt{1-\rho_k^2}$ it can be easily shown from (13) that

$$P(\sigma_{Bk}|h_k) = H\left(-\frac{\sigma_{Bk} h_k}{\lambda_k}\right) \quad (14)$$

where the H function is

$$H(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-t^2/2} dt.$$

Introducing the usual notation for the stabilities $\Delta_k \equiv \sigma_B h_k$, it follows that

$$P(\sigma_B, b_1, \{h_j\}) = P_0(b_1, h_1) P_0(h_2) P_0(h_3) \times \left[H\left(-\frac{\Delta_2}{\lambda_2}\right) H\left(-\frac{\Delta_3}{\lambda_3}\right) + \Theta(\sigma_B b_1) D\left(\frac{\Delta_2}{\lambda_2}, \frac{\Delta_3}{\lambda_3}\right) \right] \quad (15)$$

where Θ is the Heaviside function and

$$D(y_2, y_3) \equiv H(-y_2)H(y_3) + H(y_2)H(-y_3). \quad (16)$$

Integrating over b_1 we obtain

$$P(\sigma_B, \{h_j\}) = P_0(h_1) P_0(h_2) P_0(h_3) \times \left[H\left(-\frac{\Delta_2}{\lambda_2}\right) H\left(-\frac{\Delta_3}{\lambda_3}\right) + H\left(-\frac{\Delta_1}{\lambda_1}\right) D\left(\frac{\Delta_2}{\lambda_2}, \frac{\Delta_3}{\lambda_3}\right) \right]. \quad (17)$$

Introducing further simplifications in the notation, $y_k \equiv \Delta_k/\lambda_k$, we finally arrive at the result

$$P(b_1|\sigma_B, \{h_j\}) = P_0(b_1|h_1) \left[\frac{H(-y_2)H(-y_3) + \Theta(\sigma_B b_1)D(y_2, y_3)}{H(-y_2)H(-y_3) + H(-y_1)D(y_2, y_3)} \right]. \quad (18)$$

We are able now to evaluate the optimal weight, according to (12). Straightforward calculations using (13), (16) and (18) yield

$$\tilde{f}_k = \frac{3}{\sqrt{2\pi}} J_k \lambda_k \sigma_B e^{-(\Delta_k^2/2\lambda_k^2)} v_k \left(\frac{\Delta_1}{\lambda_1}, \frac{\Delta_2}{\lambda_2}, \frac{\Delta_3}{\lambda_3} \right) \quad (19)$$

where

$$v_1(y_1, y_2, y_3) \equiv \frac{D(y_2, y_3)}{H(-y_2)H(-y_3) + H(-y_1)D(y_2, y_3)}. \quad (20)$$

Note the very interesting fact that although the initial dynamics, as described by (6), did not include the traditional Hebbian term, i.e. $\sigma_B S_i$, the optimization procedure led in a natural way to its appearance, since the modulation function \tilde{F}_k has a σ_B factor. We thus arrive at a highly non-local learning rule, where information from the stabilities of the other branches is needed in order to perform the update of the synaptic couplings of a given branch. Any algorithm which does not use such non-local information would necessarily have a poorer performance.

Inserting result (19) in the mean evolution (7), we see that the optimization procedure left the ρ_1 equation uncoupled from the $\{J_k\}$ equations. Integrating over b_1 and summing over σ_B we obtain

$$\frac{d\rho_1}{d\alpha} = \frac{3}{4\pi} \rho_1 \lambda_1^2 \int \left(\prod_{j=1}^3 Dh_j \right) e^{-h_j^2/\lambda_j^2} \frac{D^2\left(\frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right)}{r\left(\frac{h_1}{\lambda_1}, \frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right) \left[1 - r\left(\frac{h_1}{\lambda_1}, \frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right)\right]} \quad (21)$$

where Dx is the Gaussian measure

$$Dx = dx \frac{e^{-x^2/2}}{\sqrt{2\pi}}$$

and

$$r\left(\frac{h_1}{\lambda_1}, \frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right) \equiv P(\sigma_B = -1 | \{h_j\}) = H\left(\frac{h_2}{\lambda_2}\right) H\left(\frac{h_3}{\lambda_3}\right) + H\left(\frac{h_1}{\lambda_1}\right) D\left(\frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right). \quad (22)$$

Inserting the optimal weight in (8), the choice $\Omega_k = 0$ leads to the equation for the evolution of the branches norms

$$\frac{dJ_1}{d\alpha} = \frac{3}{4\pi} J_1 \lambda_1^2 \int \left(\prod_{j=1}^3 Dh_j \right) e^{-h_j^2/\lambda_j^2} \frac{D^2\left(\frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right)}{r\left(\frac{h_1}{\lambda_1}, \frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right) \left[1 - r\left(\frac{h_1}{\lambda_1}, \frac{h_2}{\lambda_2}, \frac{h_3}{\lambda_3}\right)\right]} \quad (23)$$

so that, as can be seen from (21) and (23),

$$J_k(\alpha) = c \rho_k(\alpha) \quad (24)$$

where c is a constant and we can self-consistently choose $c = 1$, since ρ_k enters the J_k dynamics (23). Such a choice leaves these two equations uncoupled and, what is even more interesting, identical. This means that, in the mean,

$$J_k(\alpha) = \rho_k(\alpha) \quad (25)$$

as long as $J_k(0) = \rho_k(0)$.

4. Results and simulations

The branch permutation symmetry displayed by the set of differential equations (7) or (21) allows for a branch symmetric solution whenever the initial conditions are themselves symmetric, i.e. $\rho_k(\alpha = 0)$ is independent of k . The generalization error (5) can be obtained once the overlaps have been obtained by numerically integrating equation (21). The result is shown in figure 1 (lower full curve). The symbols represent the results of two different numerical simulations which were performed as follows. The first type of simulation

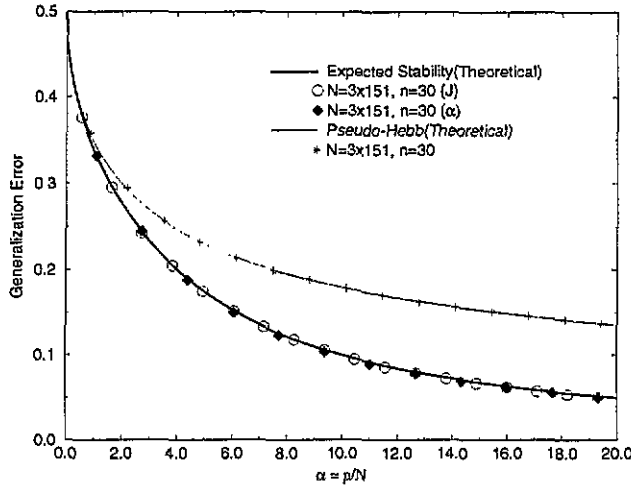


Figure 1. Theoretical results (full curves) and simulations (symbols) for the pseudo-Hebb (upper curve) and expected stability (lower curve) algorithms. N is the number of Boolean input units and n is the number of networks used for obtaining the mean ρ_k 's. Standard error bars would be approximately the same size of the symbols. See text for details.

(represented in figure 1 by diamonds) uses the value of $\rho(\alpha)$ obtained from the numerical integration of (21), in evaluating the weight function \tilde{F}_k :

$$\tilde{F}_1(\sigma_B, \tilde{h}_1, \tilde{h}_2, \tilde{h}_3, \alpha) = \frac{3}{\sqrt{2\pi}} \mathcal{D}(\alpha) \sigma_B e^{-(\tilde{h}_1^2/2\mathcal{D}^2(\alpha))} v_1 \left(\frac{\sigma_B \tilde{h}_1}{\mathcal{D}(\alpha)}, \frac{\sigma_B \tilde{h}_2}{\mathcal{D}(\alpha)}, \frac{\sigma_B \tilde{h}_3}{\mathcal{D}(\alpha)} \right) \quad (26)$$

where $\mathcal{D}(\alpha) \equiv \sqrt{1 - \rho^2(\alpha)}$ and $\tilde{h}_k \equiv J_k \cdot S_k$. The second type of simulation (represented by circles) is done using the decoupling property of the optimized dynamics and instead of using the *a priori* known value of $\rho(\alpha)$ it uses the accessible values of the norms $\{J_k\}$ as measured during run time. That means that the weight function is

$$\tilde{F}_1(\sigma_B, \tilde{h}_1, \tilde{h}_2, \tilde{h}_3, J_1, J_2, J_3) = \frac{3}{\sqrt{2\pi}} \mathcal{D}_1 \sigma_B e^{-(\tilde{h}_1^2/2\mathcal{D}_1^2)} v_1 \left(\frac{\sigma_B \tilde{h}_1}{\mathcal{D}_1}, \frac{\sigma_B \tilde{h}_2}{\mathcal{D}_2}, \frac{\sigma_B \tilde{h}_3}{\mathcal{D}_3} \right) \quad (27)$$

where $\mathcal{D}_k \equiv \sqrt{1 - J_k^2}$ and in this case the value of J_k in each step is measured and used. The fact that both simulations give rise to identical results, within error bars, shows that the last algorithm is immune to the fluctuations in the values of J_k . In this form we have a perfectly well defined and applicable algorithm.

Still in figure 1, and for the sake of comparison, we present the theoretical results accompanied by those of simulations, for what we dubbed the pseudo-Hebb algorithm. This is just a generalization of the Hebb rule to the committee machine. The Hebbian rule uses a weight function independent of the pos-synaptic fields $\{h_k\}$. If we allow for a possibly varying step size, as a function of α but not of the stabilities, the most general weight function, within what might be called the pseudo-Hebb class, is given by

$$F_k^{pH} = \sigma_B W_k(\alpha). \quad (28)$$

Insertion of weight (28) in (7) and (8) gives

$$\frac{d\rho_k}{d\alpha} = \rho_k \frac{W_k(\alpha)}{J_k} \left[\frac{1}{\sqrt{2\pi}} \left(\frac{1 - \rho_k^2}{\rho_k} \right) - \frac{W_k(\alpha)}{6J_k} \right] \quad (29)$$

$$\frac{dJ_k}{d\alpha} = \frac{W_k(\alpha)}{J_k} \left[\frac{J_k \rho_k}{\sqrt{2\pi}} + \frac{W_k(\alpha)}{6} \right]. \tag{30}$$

The function $W_k^*(\alpha)$ that optimizes $d\rho_k/d\alpha$ is

$$W_k^*(\alpha) = \frac{3}{\sqrt{2\pi}} J_k \left(\frac{1 - \rho_k^2}{\rho_k} \right) \tag{31}$$

in which case the overlaps' and the vector norms' dynamics are given by

$$\frac{d\rho_k}{d\alpha} = \frac{3}{4\pi} \frac{(1 - \rho_k^2)^2}{\rho_k} \tag{32}$$

$$\frac{dJ_k}{d\alpha} = \frac{3}{4\pi} J_k \frac{(1 - \rho_k^2)^2}{\rho_k^2} + \frac{3}{2\pi} J_k (1 - \rho_k^2). \tag{33}$$

The initial condition $\rho(0) = 0$ leads to

$$\rho_k(\alpha) = \left(1 + \frac{2\pi}{3\alpha} \right)^{-1/2}. \tag{34}$$

The weight function (31) is an awkward expression of a constant value! It is remarkable that optimization of W_k , in the sense of mean generalization, gives $W_k(\alpha) = \text{constant}$, that is, an optimal weight independent of α . An analogous result has been previously obtained in [6] for the Boolean perceptron with no hidden layers. This can be seen from (31)–(33), from which one easily obtains

$$\frac{d}{d\alpha} \ln \left(\frac{1 - \rho_k^2}{\rho_k} \right) = -\frac{d}{d\alpha} \ln(J_k) \Rightarrow \frac{d}{d\alpha} W_k^*(\alpha) = 0.$$

Result (34) and equation (5) yield an asymptotic error $e_g^{\text{pH}} \simeq \sqrt{3/2\pi} \alpha^{-1/2}$. This should be compared to the error decay of the optimized algorithm which gives, for large α ,

$$e_g \simeq \frac{\kappa_0}{\alpha} \tag{35}$$

where

$$\kappa_0 = 2 \left[\int_{-\infty}^{+\infty} Dx \frac{e^{-x^2/2}}{H(x)} \right]^{-1} \simeq 0.88. \tag{36}$$

This is *exactly* the same nominal error of the expected stability algorithm for the perceptron with no hidden layers.

The above results (equations (35) and (36)) for the asymptotical behaviour of the generalization error arise from (21) in the limit $\rho_k = \rho \rightarrow 1$ ($\lambda \rightarrow 0$). The asymptotical behaviour of the three-dimensional integral can be obtained by a single change of variables, namely $x_1 = h_1/\lambda_1$, for the first branch. The contribution for the integral of the quadrants with $h_2 h_3 < 0$ can be easily seen to be $2\lambda/\kappa_0$, while for $h_2 h_3 > 0$ the integral is of higher order in λ , thus being disregarded.

5. Conclusions

The analysis of the optimal weight function (19) yields some very interesting results. This analysis is based on the interpretation of the weight function as a measure of the information value of a new example, as judged from the previous experience of the student network acquired during the learning process. Two elements govern the behaviour of the weight function. In the first place, it depends on the error of generalization through its dependence

on the ρ_k 's. This means that the synaptic couplings will be changed differently depending on whether the net is in an advanced stage of learning or just beginning. These remarks are the basis of the construction of a self-adaptive algorithm, along the lines of [9]. Figure 2 shows the weight function dependence on the branch stability Δ_1 in the special case where $\Delta_2 = -3$ and $\Delta_3 = 3$, for different learning stages as parametrized by the overlaps ρ_k . Note the evolution of F_1 from a flat, pure pseudo-Hebbian (i.e. constant) weight at the beginning of the process to a highly structured function of the stability Δ_1 at later times. The second element is what might be dubbed the element of surprise of a new example. Consider figure 3, where the weight function (rescaled by $1/\sqrt{1-\rho_1^2}$) is shown as a function of $y_1 = \Delta_1/\lambda_1$ for various values of the pair y_2 and y_3 . By presenting the values in such a rescaled form we can analyse the effect of the surprise factor independently of the stage of learning. Notice that if a given branch variable y_k is positive that branch gives a contribution which agrees with the teacher output. In general, if the internal representation agrees with the teacher then no major change in the synaptic couplings is performed. However, and especially when the other two branches disagree (either between themselves or with the teacher), the cases with $y_1 < 0$ lead to recognizing a surprise and thus performing a major synaptic change. However, if y_1 gets to be even more negative than y_2 and y_3 , then the increased confidence of being correct, although with $\sigma_{J1} = -\sigma_B$, leads to a decrease in the weight function as if this meant that the blame for a wrong classification lies somewhere on the other branches.

A point that deserves to be mentioned is the fact that as in all the optimized on-line dynamics for spherical perceptrons [6, 8, 9], the energy function per example which characterizes the gradient descent dynamics learning algorithm is related to the conditional probability $P(\sigma_B|\{h_j\})$ by

$$E = -\lambda^2 \ln P(\sigma_B|\{h_j\}).$$

In conclusion the study of the learning dynamics of a net with three hidden units has led to algorithms that can be used in a practical learning task where this architecture is suited. The uncoupling of the ρ and J variables due to the optimization procedure rids the algorithm from unknown parameters and yields a weight function that can be evaluated with the available information. The above results were obtained for the case $K = 3$. It

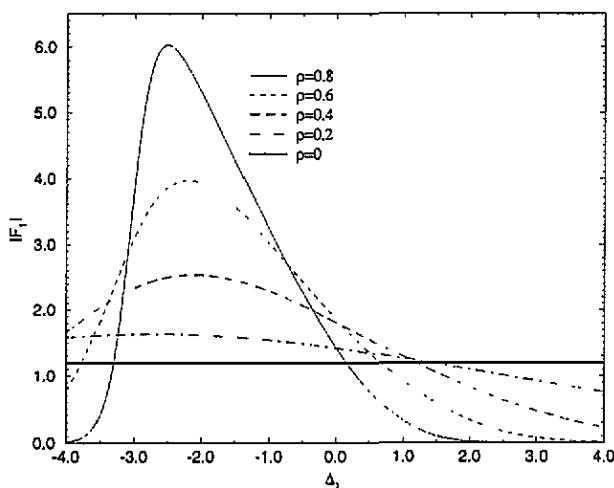


Figure 2. Optimal weight function for the first branch in the symmetrical situation where $\rho_1 = \rho_2 = \rho_3 = \rho$.

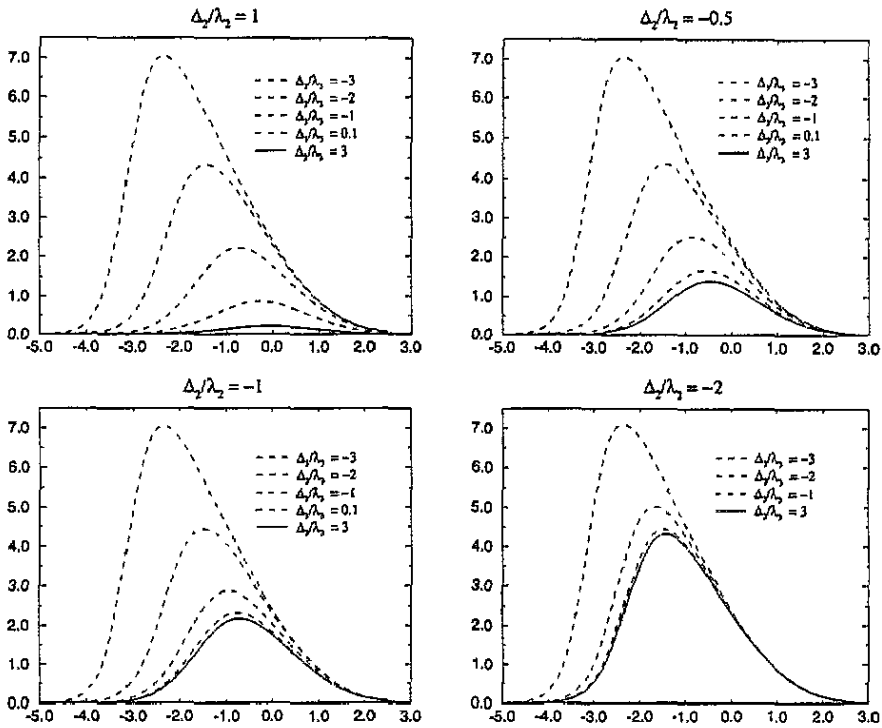


Figure 3. Rescaled optimal weight function $(1 - \rho^2)^{-1/2} |F_1|$ against rescaled stability $y_1 = \Delta_1/\lambda_1$ for several values of Δ_2/λ_2 and Δ_3/λ_3 . Full curves represent the asymptotic behaviour for increasing Δ_3/λ_3 with fixed Δ_2/λ_2 .

remains to ask to what amount they are representative of the general K machine, and how involved would their derivation be in the general case. Preliminary results [2] show that the general features of the $K = 3$ case remain true in the general case. The more general problem of two-layer networks deserves further study not only from a theoretical point of view but also due to its many possible applications. The optimization techniques we have studied here can be extended to more realistic nets and will possibly be useful due to the effectiveness and simplicity of on-line optimized algorithms. The problem of more general example distributions can be dealt with in a similar fashion and problems of catastrophic interference can probably be successfully overcome by the use of these adaptive procedures.

Acknowledgments

The authors thank O Kinouchi for several useful discussions. The work of M Copelli was supported by FAPESP, while N Caticha received support from CNPq.

References

- [1] Biehl M and Schwarze H 1994 Learning by online gradient descent *Preprint* Institut für theoretische Physik, University of Würzburg
- [2] Copelli M, Kinouchi O and Caticha N, in preparation

- [3] Gardner E and Derrida B 1989 Three unfinished works on the optimal storage capacity of networks *J. Phys. A: Math. Gen.* **22** 1983
- [4] Györgyi G and Tishby N 1989 Statistical theory of learning a rule *Neural Networks and Spin Glasses* ed W Theumann and R Köberle (Singapore: World Scientific)
- [5] Kabashima Y 1994 Perfect loss of generalization due to noise in $K = 2$ parity machines *J. Phys. A: Math. Gen.* **27** 1917
- [6] Kinouchi O 1992 Optimal generalization in perceptrons *Master's thesis* IFQSC, Universidade de São Paulo, São Carlos/SP
- [7] Kinouchi O and Caticha N 1991 Biased learning in Boolean perceptrons *Physica* **185A** 411
- [8] Kinouchi O and Caticha N 1992 Optimal generalization in perceptrons *J. Phys. A: Math. Gen.* **25** 6243
- [9] Kinouchi O and Caticha N 1993 Lower bounds on generalization errors for drifting rules *J. Phys. A: Math. Gen.* **26** 6161
- [10] Kinzel W and Ruján P 1990 Improving a network generalization ability by selecting examples *Europhys. Lett.* **13** 473
- [11] Mato G and Parga N 1992 Generalization properties of multilayered neural networks *J. Phys. A: Math. Gen.* **25** 5047
- [12] Oppen M and Haussler D 1991 Generalization performance of Bayes optimal classification algorithm for learning a perceptron *Phys. Rev. Lett.* **66** 2677