
On-Line Learning with Restricted Training Sets: Exact Solution as Benchmark for General Theories

H.C. Rae
hamish.rae@kcl.ac.uk

P. Sollich
psollich@mth.kcl.ac.uk

A.C.C. Coolen
tcoolen@mth.kcl.ac.uk

Department of Mathematics
King's College London
The Strand
London WC2R 2LS, UK

Abstract

We solve the dynamics of on-line Hebbian learning in perceptrons exactly, for the regime where the size of the training set scales linearly with the number of inputs. We consider both noiseless and noisy teachers. Our calculation cannot be extended to non-Hebbian rules, but the solution provides a nice benchmark to test more general and advanced theories for solving the dynamics of learning with restricted training sets.

1 Introduction

Considerable progress has been made in understanding the dynamics of supervised learning in layered neural networks through the application of the methods of statistical mechanics. A recent review of work in this field is contained in [1]. For the most part, such theories have concentrated on systems where the training set is much larger than the number of updates. In such circumstances the probability that a question will be repeated during the training process is negligible and it is possible to assume for large networks, via the central limit theorem, that the local field distribution is Gaussian. In this paper we consider *restricted training sets*; we suppose that the size of the training set scales linearly with N , the number of inputs. The probability that a question will reappear during the training process is no longer negligible, the assumption that the local fields have Gaussian distributions is not tenable, and it is clear that correlations will develop between the weights and the

questions in the training set as training progresses. In fact, the non-Gaussian character of the local fields should be a *prediction* of any satisfactory theory of learning with restricted training sets, as this is clearly demanded by numerical simulations. Several authors [2, 3, 4, 5, 6, 7] have discussed learning with restricted training sets but a general theory is difficult. A simple model of learning with restricted training sets which can be solved *exactly* is therefore particularly attractive and provides a yardstick against which more difficult and sophisticated general theories can, in due course, be tested and compared. We show how this can be accomplished for on-line Hebbian learning in perceptrons with restricted training sets and we obtain exact solutions for the generalisation error and the training error for a class of noisy teachers and students with arbitrary weight decay. Our theory is in excellent agreement with numerical simulations and our prediction of the probability density of the student field is a striking confirmation of them, making it clear that we are indeed dealing with local fields which are non-Gaussian.

2 Definitions

We study on-line learning in a student perceptron S , which tries to perform a task defined by a teacher perceptron characterised by a fixed weight vector $\mathbf{B}^* \in \mathbb{R}^N$. We assume, however, that the teacher is noisy and that the *actual* teacher output T and the corresponding student response S are given by

$$\begin{aligned} T : \{-1, 1\}^N &\rightarrow \{-1, 1\} & T(\boldsymbol{\xi}) &= \text{sgn}[\mathbf{B} \cdot \boldsymbol{\xi}], \\ S : \{-1, 1\}^N &\rightarrow \{-1, 1\} & S(\boldsymbol{\xi}) &= \text{sgn}[\mathbf{J} \cdot \boldsymbol{\xi}], \end{aligned}$$

where the vector \mathbf{B} is drawn *independently* of $\boldsymbol{\xi}$ with probability $p(\mathbf{B})$ which may depend explicitly on the correct teacher vector \mathbf{B}^* . Of particular interest are the following two choices, described in literature as output noise and Gaussian input noise, respectively:

$$p(\mathbf{B}) = \lambda \delta(\mathbf{B} + \mathbf{B}^*) + (1 - \lambda) \delta(\mathbf{B} - \mathbf{B}^*) \quad (1)$$

where $\lambda \geq 0$ represents the probability that the teacher output is incorrect, and

$$p(\mathbf{B}) = \left[\frac{N}{2\pi\Sigma^2} \right]^{\frac{N}{2}} e^{-\frac{N}{2}(\mathbf{B} - \mathbf{B}^*)^2/\Sigma^2}. \quad (2)$$

The variance Σ^2/N has been chosen so as to achieve appropriate scaling for $N \rightarrow \infty$.

Our learning rule will be the on-line Hebbian rule, i.e.

$$\mathbf{J}(\ell+1) = \left(1 - \frac{\gamma}{N}\right)\mathbf{J}(\ell) + \frac{\eta}{N} \boldsymbol{\xi}(\ell) \text{sgn}[\mathbf{B}(\ell) \cdot \boldsymbol{\xi}(\ell)] \quad (3)$$

where the non-negative parameters γ and η are the decay rate and the learning rate, respectively. At each iteration step ℓ an input vector $\boldsymbol{\xi}(\ell)$ is picked at random from a training set consisting of $p = \alpha N$ randomly drawn vectors $\boldsymbol{\xi}^\mu \in \{-1, 1\}^N$, $\mu = 1, \dots, p$. This set remains unchanged during the learning dynamics. At the same time the teacher selects at random, and independently of $\boldsymbol{\xi}(\ell)$, the vector $\mathbf{B}(\ell)$, according to the probability distribution $p(\mathbf{B})$. Iterating equation (3) gives

$$\mathbf{J}(m) = \left(1 - \frac{\gamma}{N}\right)^m \mathbf{J}_0 + \frac{\eta}{N} \sum_{\ell=0}^{m-1} \left(1 - \frac{\gamma}{N}\right)^{m-\ell-1} \boldsymbol{\xi}(\ell) \text{sgn}[\mathbf{B}(\ell) \cdot \boldsymbol{\xi}(\ell)] \quad (4)$$

We assume that the (noisy) teacher output is *consistent* in the sense that if a question $\boldsymbol{\xi}$ reappears at some stage during the training process the teacher makes the same choice of \mathbf{B} in both cases, i.e. if $\boldsymbol{\xi}(\ell) = \boldsymbol{\xi}(\ell')$ then also $\mathbf{B}(\ell) = \mathbf{B}(\ell')$. This consistency allows us to define a generalised training set \tilde{D} by including with the p

questions the corresponding teacher vectors:

$$\tilde{D} = \{(\boldsymbol{\xi}^1, \mathbf{B}^1), \dots, (\boldsymbol{\xi}^p, \mathbf{B}^p)\}$$

There are two sources of randomness in this problem. First of all there is the random realisation of the ‘path’ $\Omega = \{(\boldsymbol{\xi}(0), \mathbf{B}(0)), (\boldsymbol{\xi}(1), \mathbf{B}(1)), \dots, (\boldsymbol{\xi}(\ell), \mathbf{B}(\ell)), \dots\}$. This is simply the randomness of the stochastic process that gives the evolution of the vector \mathbf{J} . Averages over this process will be denoted as $\langle \dots \rangle$. Secondly there is the randomness in the composition of the training set. We will write averages over all training sets as $\langle \dots \rangle_{\text{sets}}$. We note that

$$\langle f[\boldsymbol{\xi}(\ell), \mathbf{B}(\ell)] \rangle = \frac{1}{p} \sum_{\mu=1}^p f(\boldsymbol{\xi}^\mu, \mathbf{B}^\mu) \quad (\text{for all } \ell)$$

and that averages over all possible realisations of the training set are given by

$$\begin{aligned} & \langle f[(\boldsymbol{\xi}^1, \mathbf{B}^1), (\boldsymbol{\xi}^2, \mathbf{B}^2), \dots, (\boldsymbol{\xi}^p, \mathbf{B}^p)] \rangle_{\text{sets}} \\ &= \sum_{\boldsymbol{\xi}^1} \sum_{\boldsymbol{\xi}^2} \dots \sum_{\boldsymbol{\xi}^p} \frac{1}{2^{Np}} \int \left[\prod_{\mu=1}^p p(\mathbf{B}^\mu) d\mathbf{B}^\mu \right] f[(\boldsymbol{\xi}^1, \mathbf{B}^1), (\boldsymbol{\xi}^2, \mathbf{B}^2), \dots, (\boldsymbol{\xi}^p, \mathbf{B}^p)] \end{aligned}$$

where $\boldsymbol{\xi}^\mu \in \{-1, 1\}^N$. We normalise \mathbf{B}^* so that $[\mathbf{B}^*]^2 = 1$ and choose the time unit $t = m/N$. We finally assume that \mathbf{J}_0 and \mathbf{B}^* are statistically independent of the training vectors $\boldsymbol{\xi}^\mu$, and that they obey $J_i(0), B_i^* = \mathcal{O}(N^{-\frac{1}{2}})$ for all i .

3 Explicit Microscopic Expressions

At the m -th stage of the learning process the two simple scalar observables $Q[\mathbf{J}] = \mathbf{J}^2$ and $R[\mathbf{J}] = \mathbf{B}^* \cdot \mathbf{J}$, and the joint distribution of fields $x = \mathbf{J} \cdot \boldsymbol{\xi}$, $y = \mathbf{B}^* \cdot \boldsymbol{\xi}$, $z = \mathbf{B} \cdot \boldsymbol{\xi}$ (calculated over the questions in the training set \tilde{D}), are given by

$$Q[\mathbf{J}(m)] = \mathbf{J}^2(m) \quad R[\mathbf{J}(m)] = \mathbf{B}^* \cdot \mathbf{J}(m) \quad (5)$$

$$P[x, y, z; \mathbf{J}(m)] = \frac{1}{p} \sum_{\mu=1}^p \delta[x - \mathbf{J}(m) \cdot \boldsymbol{\xi}^\mu] \delta[y - \mathbf{B}^* \cdot \boldsymbol{\xi}^\mu] \delta[z - \mathbf{B}^\mu \cdot \boldsymbol{\xi}^\mu] \quad (6)$$

For infinitely large systems one can prove that the fluctuations in mean-field observables such as $\{Q, R, P\}$, due to the randomness in the dynamics, will vanish [6]. Furthermore one assumes, with convincing support from numerical simulations, that for $N \rightarrow \infty$ the evolution of such observables, observed for different random realisations of the training set, will be reproducible (i.e. the sample-to-sample fluctuations will also vanish, which is called ‘self-averaging’). Both properties are central ingredients of all current theories. We are thus led to the introduction of the averages of the observables in (5,6), with respect to the dynamical randomness and with respect to the randomness in the training set (to be carried out in precisely this order):

$$Q(t) = \lim_{N \rightarrow \infty} \langle \langle Q[\mathbf{J}(tN)] \rangle \rangle_{\text{sets}} \quad R(t) = \lim_{N \rightarrow \infty} \langle \langle R[\mathbf{J}(tN)] \rangle \rangle_{\text{sets}} \quad (7)$$

$$P_t(x, y, z) = \lim_{N \rightarrow \infty} \langle \langle P[x, y, z; \mathbf{J}(tN)] \rangle \rangle_{\text{sets}} \quad (8)$$

A fundamental ingredient of our calculations will be the average $\langle \xi_i \operatorname{sgn}(\mathbf{B} \cdot \boldsymbol{\xi}) \rangle_{(\boldsymbol{\xi}, \mathbf{B})}$, calculated over all realisations of $(\boldsymbol{\xi}, \mathbf{B})$. We find, for a wide class of $p(\mathbf{B})$, that

$$\langle \xi_i \operatorname{sgn}(\mathbf{B} \cdot \boldsymbol{\xi}) \rangle_{(\boldsymbol{\xi}, \mathbf{B})} = \rho B_i^* + \mathcal{O}(N^{-3/2}) \quad (9)$$

where, for example,

$$\rho = \sqrt{\frac{2}{\pi}} (1-2\lambda) \quad (\text{output noise}) \quad (10)$$

$$\rho = \sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{1+\Sigma^2}} \quad (\text{Gaussian input noise}) \quad (11)$$

4 Averages of Simple Scalar Observables

Calculation of $Q(t)$ and $R(t)$ using (4, 5, 7, 9) to execute the path average and the average over sets is relatively straightforward, albeit tedious. We find that

$$\begin{aligned} Q(t) = & e^{-2\gamma t} Q_0 + 2\eta\rho R_0 \frac{e^{-\gamma t}(1-e^{-\gamma t})}{\gamma} + \frac{\eta^2}{2\gamma}(1-e^{-2\gamma t}) \\ & + \eta^2 \frac{(1-e^{-\gamma t})^2}{\gamma^2} \left(\frac{1}{\alpha} + \rho^2\right) \end{aligned} \quad (12)$$

and that

$$R(t) = e^{-\gamma t} R_0 + \eta\rho\gamma^{-1}(1-e^{-\gamma t}) \quad (13)$$

where ρ is given by equations (10, 11) in the examples of output noise and Gaussian input noise, respectively. We note that the generalisation error is given by

$$E_g = \frac{1}{\pi} \arccos \left[R(t) / \sqrt{Q(t)} \right] \quad (14)$$

All models of the teacher noise which have the same ρ will thus have the same generalisation error at any time. This is true, in particular, of output noise and Gaussian input noise when their respective parameters λ and Σ are related by

$$1-2\lambda = \frac{1}{\sqrt{1+\Sigma^2}}. \quad (15)$$

With each type of teacher noise for which (9) holds, one can thus associate an effective output noise parameter λ . Note, however, that this effective teacher error probability λ will in general not be identical to the *true* teacher error probability associated with a given $p(\mathbf{B})$, as can immediately be seen by calculating the latter for the Gaussian input noise (2).

5 Average of the Joint Field Distribution

The calculation of the average of the joint field distribution starting from equation (8) is more difficult. Writing $\sigma = (1-\gamma/N)$, and expressing the δ functions in terms of complex exponentials, we find that

$$\begin{aligned} P_t(x, y, z) = & \int \frac{d\hat{x}d\hat{y}d\hat{z}}{8\pi^3} e^{i(x\hat{x}+y\hat{y}+z\hat{z})} \lim_{N \rightarrow \infty} \left\langle e^{-i[\hat{x}e^{-\gamma t} \mathbf{J}_0 \cdot \boldsymbol{\xi}^1 + \hat{y} \mathbf{B}^* \cdot \boldsymbol{\xi}^1 + \hat{z} \mathbf{B}^1 \cdot \boldsymbol{\xi}^1]} \right. \\ & \left. \times \prod_{\ell=0}^{tN} \left[\frac{1}{p} \sum_{\nu=1}^p e^{-i[\eta\hat{x}N^{-1}\sigma^{tN-\ell}(\boldsymbol{\xi}^1 \cdot \boldsymbol{\xi}^\nu) \operatorname{sgn}(\mathbf{B}^\nu \cdot \boldsymbol{\xi}^\nu)]} \right] \right\rangle_{\text{sets}} \end{aligned} \quad (16)$$

In this expression we replace $\boldsymbol{\xi}^1$ by $\boldsymbol{\xi}$ and \mathbf{B}^1 by \mathbf{B} , and abbreviate $S = \prod_{\ell=0}^{tN} [\dots]$. Upon writing the latter product in terms of the auxiliary variables $v_\nu = (\boldsymbol{\xi}^1 \cdot \boldsymbol{\xi}^\nu) / \sqrt{N}$ and $\omega_\nu = \mathbf{B}^\nu \cdot \boldsymbol{\xi}^\nu$, we find that for large N

$$\log S \sim \chi(\hat{x} \operatorname{sgn}[\mathbf{B} \cdot \boldsymbol{\xi}], t) - \frac{i\eta\hat{x}u_1}{\gamma}(1-e^{-\gamma t}) - \frac{\eta^2\hat{x}^2u_2}{4\gamma}(1-e^{-2\gamma t}) \quad (17)$$

where u_1, u_2 are the random variables given by

$$u_1 = \frac{1}{\alpha\sqrt{N}} \sum_{\nu>1} v_\nu \operatorname{sgn}(\omega_\nu), \quad u_2 = \frac{1}{p} \sum_{\nu>1} v_\nu^2.$$

and with

$$\chi(w, t) = \frac{1}{\alpha} \int_0^t ds [e^{-[i\eta w e^{\gamma(s-t)}] - 1}] \quad (18)$$

A study of the statistics of u_1 and u_2 shows that $\lim_{N \rightarrow \infty} u_2 = 1$, and that

$$u_1 = \rho \mathbf{B}^* \cdot \boldsymbol{\xi} + \alpha^{-1/2} u \quad (N \rightarrow \infty),$$

where u is a Gaussian random variable with mean equal to zero and variance unity. On the basis of these results and equations (16, 17) we find that

$$P_t(x, y, z) = \int \frac{d\hat{x}d\hat{y}d\hat{z}}{8\pi^3} e^{i(x\hat{x}+y\hat{y}+z\hat{z})-\frac{1}{2}\hat{x}^2[Q-R^2-e^{-2\gamma t}(Q_0-R_0^2)]+\chi(\hat{x} \operatorname{sgn}[z], t)-i\hat{x}y(R-R_0e^{-\gamma t})} \\ \times \lim_{N \rightarrow \infty} \langle e^{-i[\hat{x}e^{-\gamma t} \mathbf{J}_0 \cdot \boldsymbol{\xi} + \hat{y} \mathbf{B}^* \cdot \boldsymbol{\xi} + \hat{z} \mathbf{B} \cdot \boldsymbol{\xi}]} \rangle_{(\boldsymbol{\xi}, \mathbf{B})} \quad (19)$$

where Q and R are given by the expressions (12,13) (note: $Q-R^2$ is independent of ρ , i.e. of the distribution $p(\mathbf{B})$). Let $x_0 = \mathbf{J}_0 \cdot \boldsymbol{\xi}$, $y = \mathbf{B}^* \cdot \boldsymbol{\xi}$, $z = \mathbf{B} \cdot \boldsymbol{\xi}$. We assume that, *given* y, z is independent of x_0 . This condition, which reflects in some sense the property that the teacher noise preserves the perceptron structure, is certainly satisfied for the models which we are considering and is probably true of all reasonable noise models. The joint probability density then has the form $p(x_0, y, z) = p(x_0|y)p(y, z)$. Equation (19) then leads to the following expression for the conditional probability of x , given y and z :

$$P_t(x|y, z) = \int \frac{d\hat{x}}{2\pi} e^{i\hat{x}[x-Ry]-\frac{1}{2}\hat{x}^2[Q-R^2]+\chi(\hat{x} \operatorname{sgn}[z], t)} \quad (20)$$

We observe that this probability distribution is the same for all models with the same ρ and that the dependence on z is through $\tau = \operatorname{sgn}[z]$, a directly observable quantity. The training error and the student field probability density are given by

$$E_{\text{tr}} = \int dx dy \sum_{\tau=\pm 1} \theta(-x\tau) P_t(x|y, \tau) P(\tau|y) P(y) \quad (21)$$

$$P_t(x) = \int dy \sum_{\tau=\pm 1} P_t(x|y, \tau) P(\tau|y) P(y) \quad (22)$$

in which $P(y) = (2\pi)^{-\frac{1}{2}} e^{-\frac{1}{2}y^2}$. We note that the dependence of E_{tr} and $P_t(x)$ on the specific noise model arises solely through $P(\tau|y)$ which we find is given by

$$P(\tau|y) = \lambda\theta(-\tau y) + (1-\lambda)\theta(\tau y) \quad P(\tau|y) = \frac{1}{2}(1 + \tau \operatorname{erf}[y/\sqrt{2}\Sigma])$$

in the output noise and Gaussian input noise models, respectively. In order to simplify the numerical computation of the remaining integrals one can further reduce the number of integrations analytically. Details will be reported elsewhere.

6 Comparison with Numerical Simulations

It will be clear that there is a large number of parameters that one could vary in order to generate different simulation experiments with which to test our theory. Here we have to restrict ourselves to presenting a number of representative results. Figure 1 shows, for the output noise model, how the probability density $P_t(x)$ of

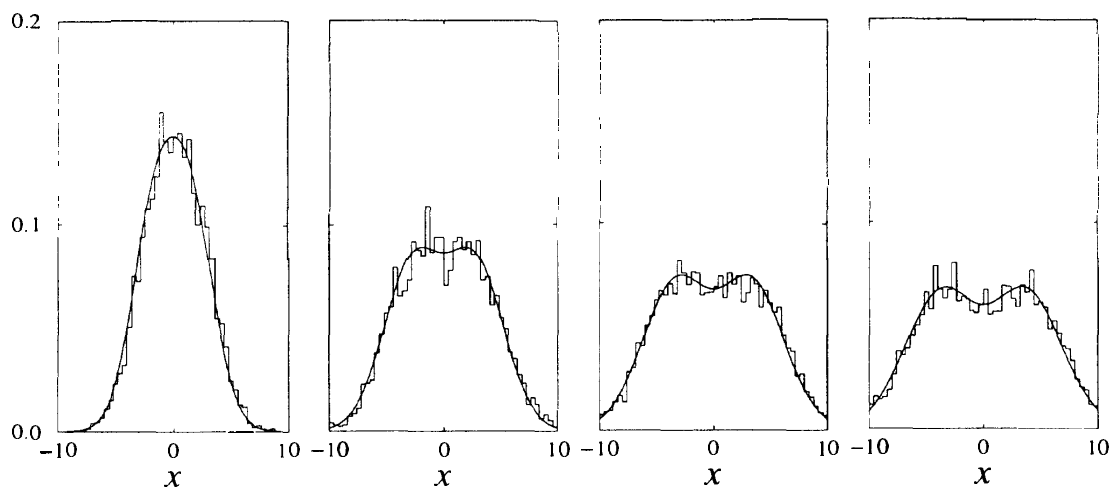


Figure 1: Student field distribution $P(x)$ for the case of output noise, at different times (left to right: $t = 1, 2, 3, 4$), for $\alpha = \gamma = \frac{1}{2}$, $J_0 = \eta = 1$, $\lambda = 0.2$. Histograms: distributions measured in simulations, ($N = 10,000$). Lines: theoretical predictions.

the student field $x = \mathbf{J} \cdot \boldsymbol{\xi}$ develops in time, starting as a Gaussian at $t = 0$ and evolving to a highly non-Gaussian distribution with a double peak by time $t = 4$. The theoretical results give an extremely satisfactory account of the numerical simulations. Figure 2 compares our predictions for the generalisation and training errors E_g and E_{tr} with the results of numerical simulations, for different initial conditions, $E_g(0) = 0$ and $E_g(0) = 0.5$, and for different choices of the two most important parameters λ (which controls the amount of teacher noise) and α (which measures the relative size of the training set). The theoretical results are again in excellent agreement with the simulations. The system is found to have no memory of its past (which will be different for some other learning rules), the asymptotic values of E_g and E_{tr} being independent of the initial student vector. In our examples E_g is consistently larger than E_{tr} , the difference becoming less pronounced as α increases. Note, however, that in some circumstances E_{tr} can also be larger than E_g . Careful inspection shows that for Hebbian learning there are no true overfitting effects, not even in the case of large λ and small γ (for large amounts of teacher noise, without regularisation via weight decay). Minor finite time minima of the generalisation error are only found for very short times ($t < 1$), in combination with special choices for parameters and initial conditions.

7 Discussion

Starting from a microscopic description of Hebbian on-line learning in perceptrons with restricted training sets, of size $p = \alpha N$ where N is the number of inputs, we have developed an exact theory in terms of macroscopic observables which has enabled us to predict the generalisation error and the training error, as well as the probability density of the student local fields in the limit $N \rightarrow \infty$. Our results are in excellent agreement with numerical simulations (as carried out for systems of size $N = 5,000$) in the case of output noise; our predictions for the Gaussian input noise model are currently being compared with the results of simulations. Generalisations of our calculations to scenarios involving, for instance, time-dependent learning rates or time-dependent decay rates are straightforward. Although it will be clear that our present calculations cannot be extended to non-Hebbian rules, since they

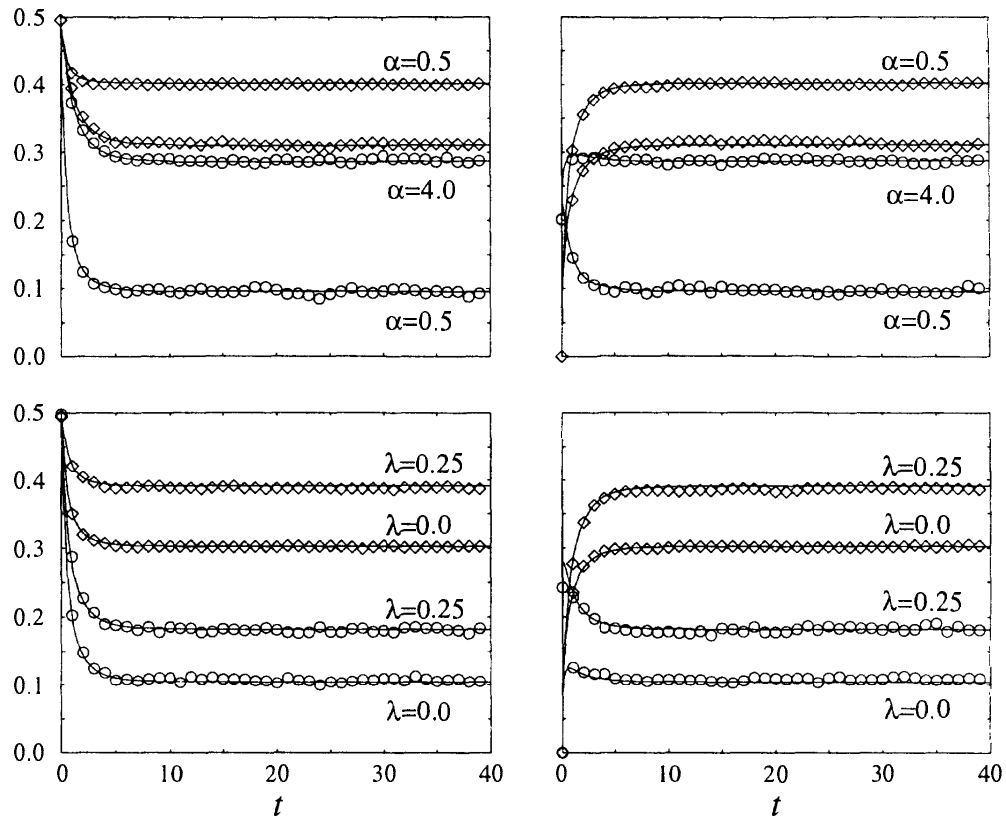


Figure 2: Generalisation errors (diamonds/lines) and training errors (circles/lines) as observed during on-line Hebbian learning, as functions of time. Upper two graphs: $\lambda = 0.2$ and $\alpha \in \{0.5, 4.0\}$ (upper left: $E_g(0) = 0.5$, upper right: $E_g(0) = 0$). Lower two graphs: $\alpha = 1$ and $\lambda \in \{0.0, 0.25\}$ (lower left: $E_g(0) = 0.5$, lower right: $E_g(0) = 0$). Markers: simulation results for an $N = 5,000$ system. Solid lines: predictions of the theory. In all cases $J_0 = \eta = 1$ and $\gamma = 0.5$.

ultimately rely on our ability to write down the microscopic weight vector \mathbf{J} at any time in explicit form (4), they do indeed provide a significant yardstick against which more sophisticated and more general theories can be tested. In particular, they have already played a valuable role in assessing the conditions under which a recent general theory of learning with restricted training sets, based on a dynamical version of the replica formalism, is exact [6, 7].

References

- [1] Mace C.W.H. and Coolen A.C.C. (1998) *Statistics and Computing* **8**, 55
- [2] Horner H. (1992a), *Z.Phys. B* **86**, 291; (1992b), *Z.Phys. B* **87**, 371
- [3] Krogh A. and Hertz J.A. (1992) *J.Phys. A: Math. Gen.* **25**, 1135
- [4] Sollich P. and Barber D. (1997) *Europhys. Lett.* **38**, 477
- [5] Sollich P. and Barber D. (1998) *Advances in Neural Information Processing Systems 10*, Eds. Jordan M., Kearns M. and Solla S. (Cambridge: MIT)
- [6] Coolen A.C.C. and Saad D., King's College London preprint KCL-MTH-98-08
- [7] Coolen A.C.C. and Saad D. (1998) (in preparation)