

Convergence theorems for a class of learning algorithms with VLRPs ¹

J.F. Feng ^{a,b,*}, B. Tirozzi ^c

^a *Mathematisches Institut, Universität München, D-80333 München, Germany*

^b *Laboratory of Biomathematics, The Babraham Institute, Cambridge CB2 4AT, United Kingdom*

^c *Mathematical Department, University of Rome "La Sapienza", P. le A. Moro, 00185 Rome, Italy*

Received 23 January 1995; accepted 21 July 1996

Abstract

We first consider the convergence of the simple competitive learning with vanishing learning rate parameters (VLRPs). Examples show that even in this setting the learning fails to converge in general. This brings us to consider the following problem: to find out a family of VLRPs such that an algorithm with the VLRPs reaches the global minima with probability one. Here, we present an approach different from stochastic approximation theory and determine a new family of VLRPs such that the corresponding learning algorithm gets out of the metastable states with probability one. In the literature it is generally believed that a family of reasonable VLRPs is of the order of $1/t^\alpha$ for $1/2 < \alpha \leq 1$, where t is the time. However, we find that a family of VLRPs which makes the algorithm go to the global minima should be between $1/\log t$ and $1/\sqrt{\log t}$.

Keywords: Vanishing learning rate parameters (VLRPs); Simple competitive learning; Simulated annealing; Stochastic approximation theory

1. Introduction

In recent years, there are extensive research works devoted to the study of the self-organizing Kohonen algorithm, both theoretically and numerically [3,4,17]. In [11,12] and references given therein, the authors consider the equilibrium states of the self-organizing Kohonen algorithm with the learning rate parameter independent of time.

* Corresponding author. Email: jf218@cam.ac.uk

¹ Partially supported by the CNR of Italy, J. Feng was also partially supported by the A. von Humboldt Foundation of Germany.

In [11], a thorough investigation of the existence and the number of the metastable states is carried out. In [7–9] the asymptotic property of the one-dimensional self-organizing Kohonen algorithm is studied. Recently, a novel approach [22] to the problem of constructing topology preserving maps is introduced which is based on a Hebbian adaption rule with winner-take-all like competition. Here, we first consider the convergence problem of the simplest version of self-organizing Kohonen algorithm (see [17], pp. 217–222), the simple competitive learning with the *nonincreasing* vanishing learning rate parameters (VLRPs) $\eta(t) > 0$ satisfying the usual restrictions found in stochastic approximation theory [6,16,18,19,23,24]:

$$\int_0^{\infty} \eta(u) du = \infty, \quad (\text{I})$$

$$\int_0^{\infty} \eta^2(u) du < \infty. \quad (\text{II})$$

The constraints (I) and (II) above are usually imposed on the stochastic approximation algorithm (see, for example, [6,18,19,23]) and the reason for a family of learning rate parameters satisfying them is fully explained in [6,23]. See also Section 3 of the present paper, where we assert that the condition (II) is not a necessary one. Examples in Section 3 show that even in this simple setting, there are already metastable states for the algorithm. Note that a canonical candidate of $\eta(t)$ under the restrictions (I) and (II) will be

$$\eta(t) = \frac{1}{t^\alpha},$$

for $1/2 < \alpha \leq 1$ ([17] p. 223, and [24] p. 259).

The above conclusions naturally suggest to us to ask the question: does here exist a general rule (VLRPs) for the learning algorithm which allows the system to get out of the metastable states? In other words, we look for a family of VLRPs which has a role like the decreasing “temperature” in simulated annealing. Nevertheless, an example in Section 3 of the present paper indicates that under the constraints (I) and (II), the learning algorithm will stay at some local minima with a positive probability.

It was first noted in [24], at p. 259, that in a linear learning algorithm with VLRPs the restriction (II) above is unnecessary and it could be replaced by a much weaker condition

$$\lim_{t \rightarrow \infty} \eta(t) = 0. \quad (\text{II}')$$

Based upon the self-similarity property of Brownian motion and results of simulated annealing in [2], we present a novel and rigorous approach to determine a new family of VLRPs. This new family of VLRPs, which is between $1/\log t$ and $1/\sqrt{\log t}$, ensures that the learning algorithm with the VLRPs escapes from the local minima and reaches the desired global minima. This fact is shown in Section 3. Note that this family of VLRPs fulfils the restriction (I) and violates the restriction (II), but it satisfies (I) and (II'). We believe that our discovery is of general guidance for a class of learning

algorithms with VLRLPs, such as the learning algorithm of Oja’s law [17], self-organizing Kohonen algorithm, Hebb learning [13] and the em and EM algorithms [1] etc.

2. A convergence theorem

In this section we consider the convergence of simple competitive learning.

2.1. Notation and results

For a concrete description of our result, we first review simple competitive learning in details.

In simple competitive learning networks there is a single layer of output units $O_i(n) \in \{1,0\}$, $i = 1, \dots, N$ at time n , each fully connected to a set of inputs $\xi_j(n)$, $j = 1, \dots, M$, via connections $w_{ij}(n)$. In the sequel, we assume that the inputs $\xi_j(n)$, $j = 1, \dots, M$, are chosen independently according to a probability distribution P . Only one of the output units, called the winner, can fire at a time. The winner is the output unit with the smallest distance between its connections and the inputs

$$|w_i(n) - \xi(n)|$$

for vectors $w_i(n) = (w_{ij}(n), j = 1, \dots, M)$, $\xi(n) = (\xi_j(n), j = 1, \dots, M)$, where $|\cdot|$ represents the Euclidean norm. Let $\bar{I}(\cdot, \cdot)$ be a function given by

$$\bar{I}(w_i(n), \xi(n+1)) = I_{\{|w_i(n) - \xi(n+1)| < |w_j(n) - \xi(n+1)|, j \neq i\}}(w_i(n), \xi(n+1)),$$

where I_A is the indicator function, i.e. $I_A(x) = 1$ if $x \in A$ and $I_A(x) = 0$ if $x \notin A$.

The simple competitive learning rule ([7], pp. 217–222) is such that only the winner changes its weights

$$w_{ij}(n+1) = w_{ij}(n) + \eta(n) \bar{I}(w_i(n), \xi(n+1)) (\xi_j(n+1) - w_{ij}(n)),$$

$$i = 1, \dots, N, j = 1, \dots, M. \tag{1}$$

$\eta(n)$ is the positive learning parameter, $\eta(0) < 1$, $\eta(n) \geq \eta(n+1)$. After the learning procedure is finished any set of input vectors will be partitioned into nonoverlapping clusters. This means that a new incoming signal $\xi(n+1)$ is classified as the pattern i if it is closest to the weight w_i . In other words the new signal $\xi(n+1)$ is recognized to be of the type w_i if $\bar{I}(w_i(n), \xi(n+1)) = 1$ and this happens if and only if

$$|w_i - \xi(n+1)| \leq |w_j - \xi(n+1)|, \quad j \neq i.$$

Note that the non-linearity of the dynamics above is addressed by the function \bar{I} . In the case considered in [24], at p. 279, $N = 1$ and so the dynamics defined by Eq. (1) is linear because there is no competition at all.

For a compact region Ω of \mathbb{R}^M , let us introduce the definition of Voronoi tessellation associated with a family of vectors $y = (y_i, i = 1, \dots, N) \in \Omega$.

Definition 1. For a given compact subset $\Omega \in \mathbb{R}^M$, the Voronoi tessellation $\Pi(y) = (\Pi(y)_i, i = 1, \dots, N)$ associated with a family of vectors y_1, \dots, y_N is a partition of Ω given by

$$\Pi(y)_i = \{x, |y_i - x| \leq |y_j - x|, j \neq i\}, \quad i = 1, \dots, N.$$

Let us define a function g which is the leading term of the supermartingale difference introduced in Appendix A:

$$g(y_1, y_2, \dots, y_N; w_1, w_2, \dots, w_N) = \sum_{i=1}^N (y_i - w_i) \cdot \int_{\Pi(y)_i} (x - y_i) f(x) dx.$$

g depends on the vectors $w = (w_1, w_2, \dots, w_N)$, $y = (y_1, y_2, \dots, y_N) \in \mathbb{R}^{M \times N}$. f is the density of the distribution P with support on a compact region Ω of \mathbb{R}^M , $\Pi(y) = (\Pi(y)_i, i = 1, \dots, N)$ is the Voronoi tessellation associated with $y = (y_1, \dots, y_N)$.

We define also:

$$\Theta := \{\text{the set of all Voronoi tessellations associated with } \{w_1(n), \dots, w_N(n)\}, \text{ for all } n\}.$$

For $y_1, \dots, y_N \in \mathbb{R}^M$, we use the convention that $y = (y_1, \dots, y_N) \in \Theta$ implies that there exists a Voronoi tessellation $\Pi(y)$ such that $\{\Pi(y)_i, i = 1, \dots, N\} \in \Theta$.

Now we state the main theorem of this section.

Theorem 2. *If there exists a unique point $(w_1, w_2, \dots, w_N) \in \mathbb{R}^{M \times N}$, such that*

$$g(y_1, \dots, y_N; w_1, \dots, w_N) \leq 0, \quad \forall y \in \Theta, \quad (2)$$

where the equality holds if and only if $y_i = w_i, i = 1, \dots, N$ and

$$\sum_n \eta(n) = \infty \quad (3a), \quad \sum_n \eta^2(n) < \infty \quad (3b) \quad (3)$$

then we almost surely have

$$\lim_{n \rightarrow \infty} w_i(n) = w_i, \quad i = 1, \dots, N.$$

Proof. The proof of Theorem 2 is postponed to Appendix A.

Let us say a few words concerning the condition (2). The fulfilment of the condition (2) ensures that the learning algorithm moves downhill in the energy landscape, and so the uniqueness of the limit of the learning algorithm is true under the condition (2). In Section 3, we will give a new family of learning rate parameters when the condition (2) is violated, which is of course the more interesting case.

For a one-dimensional input signal, i.e. $M = 1$, without loss of generality, we can assume that $a \leq w_1(0) < w_2(0) < \dots < w_N(0) \leq b$ with $\Omega = [a, b]$. In this setting, we are able to simplify the condition (2) in Theorem 2 due to the fact that the simple competitive learning does not change the order of weights $a \leq w_1(n) < w_2(n) < \dots < w_N(n) \leq b, n \geq 1$.

Lemma 3. *If $M = 1$, then $y_1, \dots, y_N \in \Theta$ if and only if $a \leq y_1 < \dots < y_N \leq b$.*

Proof. The proof is postponed to Appendix A.

By combining Lemma 3 and Theorem 2, we have the following corollary. Three examples which explain the application of the next corollary are presented in Section 2.2.

Corollary 4. *If there exists a unique point $(w_1, \dots, w_N) \in \mathbb{R}^N$, such that the inequality*

$$\begin{aligned}
 &g(y_1, \dots, y_N; w_1, \dots, w_N) \\
 &= (y_1 - w_1) \int_a^{(y_1+y_2)/2} (x - y_1) f(x) dx + \sum_{i=2}^{N-1} (y_i - w_i) \\
 &\quad \times \int_{(y_{i-1}+y_i)/2}^{(y_i+y_{i+1})/2} (x - y_i) f(x) dx + (y_N - w_N) \int_{(y_N+y_{N-1})/2}^b (x - y_N) f(x) dx \\
 &< 0
 \end{aligned}$$

holds for $a \leq y_1 < \dots < y_N \leq b$ except for $y_i = w_i$, $i = 1, \dots, N$ and

$$\sum_n \eta(n) = \infty, \quad \sum_n \eta^2(n) < \infty,$$

then we almost surely have

$$\lim_{n \rightarrow \infty} w_i(n) = w_i, \quad i = 1, \dots, N.$$

In the next theorem, we consider the convergence rate of the simple competitive learning. We prove that, under the conditions in Theorem 2, the algorithm will achieve the given accuracy within a finite number of updates.

Define

$$\tau(\epsilon) = \inf\{n, |w_i(n) - w_i| \leq \epsilon, i = 1, \dots, N\}$$

as the first time that the training error is less than ϵ .

Theorem 5. *In the circumstances of Theorem 2, there exists a constant*

$$B(\epsilon) > 0,$$

such that we have almost surely

$$\tau(\epsilon) < B(\epsilon).$$

Proof. The proof of Theorem 5 is also postponed to Appendix A.

Although $w_i(n)$, $i = 1, \dots, N$ is a stochastic process Theorem 5 asserts that within a finite and a deterministic time $B(\epsilon)$ $w_i(n)$, $i = 1, \dots, N$ will reach a given accuracy ϵ .

2.2. Examples

In this section, in order to show the applications of the theorems of the previous section, we consider three typical examples in the sense that the first example takes into account the case when the input data set is discrete, the second and the third example consider the case when the input data set is continuously distributed according to the uniform distribution and the normal distribution, respectively.

Example 1. Suppose that $f(x) = \sum_{i=1}^N c_i \delta w_i(x)$, with $\sum_{i=1}^N c_i = 1$ for $w_i \in [a, b] \subset \mathbb{R}^1$, $c_i > 0$, $i = 1, \dots, N$ and $w_1 < w_2 < \dots < w_N$. Then we have that

$$\begin{aligned} g(y_1, \dots, y_N; w_1, \dots, w_N) &= -c_1 (y_1 - w_1)^2 I_{[a, (y_1 + y_2)/2]}(w_1) \\ &\quad - \sum_{i=2}^{N-1} c_i (y_i - w_i)^2 I_{[(y_{i-1} + y_i)/2, (y_i + y_{i+1})/2]}(w_i) \\ &\quad - c_N (y_N - w_N)^2 I_{[(y_N + y_{N+1})/2, b]}(w_N). \end{aligned}$$

From the theorems of the previous section we can conclude that

$$w = (w_1, \dots, w_N)$$

is the unique attracting point of the dynamics (1). The proof of Theorem 2, see Appendix A, shows that the function g is the main contribution to the derivative of a Lyapunov function. In fact the quantity

$$\sum_{i=1}^N E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n)$$

introduced in the proof can be considered to be the Lyapunov function of the system. The difference appearing in the submartingale condition:

$$\sum_{i=1}^N [E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n) - |w_i(n) - w_i|^2]$$

can be considered as a discretized derivative and is the sum of two terms. The one different from g vanishes. The points (y_1, \dots, y_N) which make the function g equal to zero can be interpreted as the minima of this Lyapunov function. Using this terminology one may say that the dynamics (1) will converge to the global minima $y_1 = w_1, \dots, y_N = w_N$ if the hypothesis of the Theorem 2 is satisfied and if there are many points for which the equality $g = 0$ is verified, then they may be seen as local minima which can trap the dynamics. The condition ensuring that there is a unique solution (y_1, \dots, y_N) of the equation

$$g(y_1, \dots, y_N; w_1, \dots, w_N) = 0$$

is quite restrictive. In general, there are (infinitely) many solutions of it. Hence, the development of an algorithm to avoid the metastable states is of general importance, which is the content of the next section.

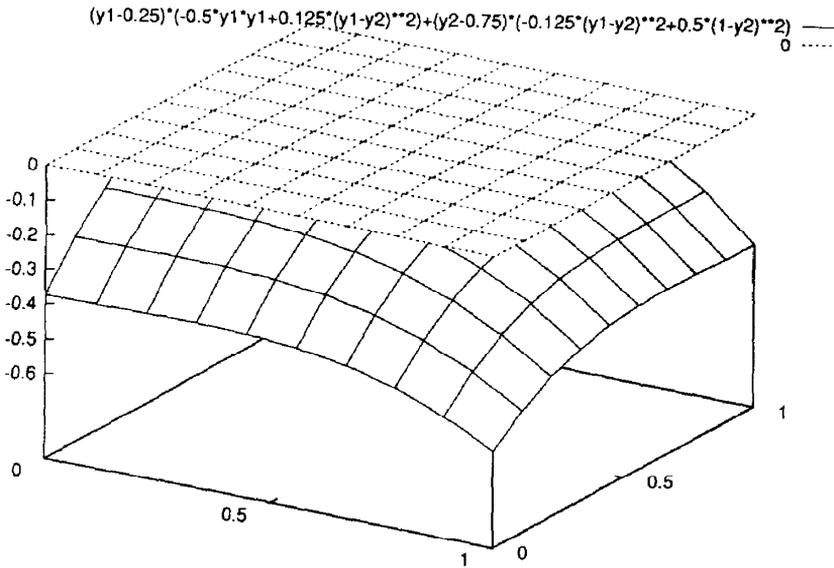


Fig. 1. The function g defined in the Example 2 of Section 2.

Example 2. Suppose that $\xi(n)$ is uniformly distributed over the interval $[0,1]$. We are going to prove that $w_1 = 1/4$, $w_2 = 3/4$ and $g(y,w)$ is negative except for $y_1 = w_1$ and $y_2 = w_2$.

First note that in this situation we have

$$g(y,w) = (y_1 - w_1) \int_0^{(y_1+y_2)/2} (x - y_1) dx + (y_2 - w_2) \int_{(y_1+y_2)/2}^1 (x - y_2) dx.$$

Therefore,

$$\begin{aligned} g(y,w) &= (y_1 - w_1) \left(\int_0^{y_1} + \int_{y_1}^{(y_1+y_2)/2} \right) (x - y_1) dx \\ &\quad + (y_2 - w_2) \left(\int_{(y_1+y_2)/2}^{y_2} + \int_{y_2}^1 \right) (x - y_2) dx \\ &= (y_1 - w_1) \left(-\frac{1}{2} y_1^2 + \frac{1}{2} \frac{(y_1 - y_2)^2}{4} \right) \\ &\quad + (y_2 - w_2) \left(-\frac{1}{2} \frac{(y_1 - y_2)^2}{4} + \frac{1}{2} (1 - y_2)^2 \right). \end{aligned}$$

It is easy to check numerically (Fig. 1) that $w_1 = 1/4$, $w_2 = 3/4$ is the unique point for $g(y,w) = 0$. Therefore, from Corollary 4 and Theorem 5 of the previous section, we have

$$\lim_{n \rightarrow \infty} w_1(n) = \frac{1}{4}, \quad \lim_{n \rightarrow \infty} w_2(n) = \frac{3}{4}$$

and $\forall \epsilon > 0, \exists B(\epsilon) > 0,$

$$\tau(\epsilon) < B(\epsilon).$$

Example 3. Suppose that $\xi(n)$ is distributed with density function

$$f(x) = \frac{1}{c} \exp\left(-\frac{x^2}{2}\right) I_{[-K, K]}(x),$$

the restriction of the normal distribution with mean 0 and variance 1 to $[-K, K]$, where $K=2$ and $c = \int_{-K}^K e^{-(x^2/2)} dx$. It is natural to expect that $w_1 = (-(1/2c))(1 - e^{-(K^2/2)}) = 0.18$ and $w_2 = (1/2c)(1 - e^{-(K^2/2)})$.

Let

$$\begin{aligned} \frac{g(y, w)}{c} &= (y_1 - w_1) \int_{-K}^{(y_1+y_2)/2} (x - y_1) e^{-(x^2/2)} dx \\ &\quad + (y_2 + w_2) \int_{(y_1+y_2)/2}^K (x - y_2) e^{-(x^2/2)} dx \\ &= -(y_1 - w_1) \frac{1}{2} e^{-((y_1+y_2)^2/8)} - (y_1 - w_1) y_1 \\ &\quad \times \int_{-K}^{(y_1+y_2)/2} e^{-(x^2/2)} dx + \frac{1}{2} e^{-(K^2/2)} (y_1 - y_2 - w_1 + w_2) \\ &\quad + (y_2 - w_2) \frac{1}{2} e^{-((y_1+y_2)^2/8)} - (y_2 - w_2) y_2 \int_{(y_1+y_2)/2}^K e^{-(x^2/2)} dx. \end{aligned}$$

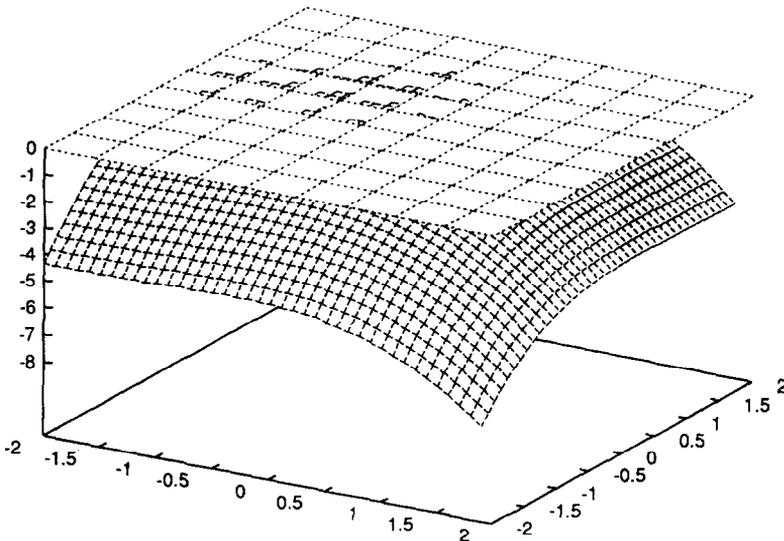


Fig. 2. The function g defined in the Example 3 of Section 2. Note that there are several points (y_1, y_2) with the property $g(y_1, y_2; w_1, w_2) = 0$.

It is easy to check numerically (see Fig. 2) that the condition of Corollary 4 on the function g is not true, i.e. there are several points (y_1, y_2) , $y_1 < y_2$ such that $g(y_1, y_2; w_1, w_2) = 0$.

3. A new family of VLRLPs

In Section 2 we developed a condition for the convergence of simple competitive learning with VLRLPs. However, it is readily seen from our examples that except for some special case (Example 2), the convergence of the algorithm will fail in general. On the other hand, all the algorithms similar to the simple competitive learning with VLRLPs are in danger of getting caught in some local minima which are useless [11]. Hence, the problem of getting out of local minima is of general importance for the learning algorithms with VLRLPs.

Essentially a learning algorithm as we considered in Section 2 with VLRLPs can be written as

$$dX_t = \eta(t)(b(X_t) dt + \beta(t) dB_t), \tag{4}$$

for $X_t \in \mathbb{R}^{M \times N}$, $b(\cdot)$ a measurable function on $\mathbb{R}^{M \times N}$, $t \in \mathbb{R}^+$ and $\eta(t)$, the VLRLPs with $\eta(0) \geq 0$, $\eta(t) \leq \eta(s)$ if $t \geq s$, $\beta(t) > 0$. Note that the discretized equation corresponding to Eq. (4) is

$$X_{n+1} = X_n + h\eta(n)b(X_n) + \beta(n)\eta(n)\sqrt{h}W_n, \tag{5}$$

where h is the step size, W_n is normally distributed with zero mean and covariance equal to the unit matrix I . For the purpose of finding a family of appropriate VLRLPs in the self-organizing Kohonen algorithm, Eq. (4) has been discussed in [24] from the Fokker–Planck equation point of view. In fact in the field of neural networks, there are many learning algorithms developed with VLRLPs and they are special cases of Eq. (4), for example, the network with Oja’s rule [17], self-organizing Kohonen algorithm, the algorithm proposed in [1] at p. 64, the dynamic link network [20,21], etc.

In this section, we first consider how to choose $\eta(t)$ and $\beta(t)$ so that X_t converges to the global minima of U if $b = -\text{grad } U$. It is proved in Theorem 6 below that under the usual restriction (I) of Section 1 on $\eta(t), \beta(t)$ should take the form (see Theorem 6)

$$\beta(t) \sim \frac{1}{\sqrt{\eta(t) \log \int_0^t \eta(u) du}}.$$

Note that as $\eta(t) = c$ a constant independent of time, Theorem 6 reduces to the well-known results of simulated annealing [2].

Secondly, if the signal is not separable from the noise, this means that in the Eq. (4) we require $\beta(t) = 1, \forall t$. It is shown in Theorem 4 below that if the family of VLRLPs $\eta(t)$ satisfies an ODE the solution of which is between $1/\log t$ and $1/\sqrt{\log t}$, then X_t will converge to the global minima with probability one. It is worthwhile to point out here that this family of VLRLPs does not satisfy the restriction (II) of Section 1 already. We believe that the discovery of this section is of general importance also for some

well-known statistical algorithms such as Robbins–Moro procedure and Kiefer–Wolfowitz procedure which have been intensively studied in the statistics ([6,16,19,23]) and take the form of Eq. (4). For the neural network applications of these algorithms we refer the reader to [5], Chapter 2.

3.1. The general case

In this section, we consider the Eq. (4):

$$dX_t = \eta(t)(b(X_t) dt + \beta(t) dB_t).$$

In order to develop a new learning rate for ensuring the convergence of the algorithm to the global minima we apply the results of simulated annealing to our case [2]. However, simulated annealing corresponds to the case in which the dynamics without noise is homogeneous, namely $\eta(t)$ is a constant independent of time t , and the noise goes to zero as the system evolves. This requires that in Eq. (4) the coefficient in front of b should be independent of t , while there is still a vanishing rate before the Brownian motion B_t . Fortunately, after taking another time scaling, we are able to remove the vanishing term in front of the drift term b , and keep the second term of the noise as a standard Brownian motion because of the self-similarity property of the Brownian motion. Furthermore, there is again a vanishing rate multiplying the Brownian motion.

Before going to more general cases, we show here first an example in order to explain our general ideas above.

Example 4. Take $\eta(t) = 1/t$, $\beta(t) = 1$, $M = N = 1$ in Eq. (4). Note that in this setting, the conditions (I) and (II) of Section 1 are fulfilled for the choice of $\eta(t)$. Now the dynamics (4) reads

$$dX_t = \frac{1}{t}(b(X_t) dt + dB_t).$$

In order to change the time scaling of the above dynamics, let

$$s = \log t = \int_0^t \eta(u) du \quad (6)$$

or

$$t = \exp(s)$$

and $Y_s = X_{e^s}$. Then

$$dX_{e^s} = dY_s = \frac{1}{e^s} b(Y_s) e^s ds + \frac{1}{e^s} dB_{e^s}. \quad (7)$$

From the self-similarity property of the Brownian motion, we know that

$$e^{-\frac{s}{2}} dB_{e^s} \sim N(0, ds).$$

So we introduce a new time scaling s and write $d\tilde{B}_s = e^{-(s/2)} dB_{e^s}$, \tilde{B}_s is again a standard Brownian motion. Now Eq. (3) can be rewritten as

$$dY_s = b(Y_s) ds + e^{-(s/2)} d\tilde{B}_s. \quad (8)$$

The relation (6) between the time t and s tells us that if s goes to infinity, then t goes to infinity also and vice versa. So if we know the limit behavior of Y_s , we know the limit behavior of X_t as well. From the general results of simulated annealing [2,14,15], we know that in the case of Eq. (8), Y_s will have positive probability to stay at any local minimum since the noise vanishes too fast, at a rate of $\exp(-(s/2))$. In order to ensure that X_t is not trapped in some local minima, we should slow down the decreasing rate of the noise. For this example, a correct choice is (see Theorem 6):

$$dZ_t = \eta(t) \left[b(Z_t) dt + \frac{\gamma}{\sqrt{\eta(t) \log \log t}} dB_t \right] \tag{9}$$

for a constant γ which as in simulated annealing is problem dependent.

In [23] and [19], under the restriction of

$$\int_0^\infty \eta(u) du = \infty, \quad \int_0^\infty \eta(u)^2 \beta(u)^2 du < \infty \tag{10}$$

for the stochastic differential Eq. ((I) and (II) of Section 1 is a special case of Eq. (10)),

$$dX_t = \eta(t) b(X_t) dt + \eta(t) \beta(t) dB_t,$$

the convergence of the solution to the set of attractors (not global minima!) of the above dynamics is proved. However, we note that in Eq. (9), the VLRPs

$$\eta(t) = \frac{1}{t}, \quad \beta(t) \eta(t) = \frac{1}{\sqrt{t \log \log t}}$$

with

$$\int_0^\infty \eta(u) du = \infty, \quad \int_0^\infty \eta(u)^2 \beta(u)^2 du = \int_0^\infty \frac{1}{u \log \log u} du = \infty$$

already violate the usual restriction (10) found in stochastic approximation theory.

In general, we have the following result for $b(x) = -\text{grad } U(x)$ for a function U defined on Ω (see Remark 2).

Theorem 6. *Suppose that*

$$\lim_{t \rightarrow \infty} \int_0^t \eta(u) du = \infty \tag{11}$$

and

$$dZ_t = \eta(t) \left[b(Z_t) dt + \frac{\gamma}{\sqrt{\eta(t) \log \int_0^t \eta(u) du}} dB_t \right], \tag{12}$$

where $Z_t \in \Omega$, a compact subset of $\mathbb{R}^{M \times N}$, b is a measurable function on $C^1(\Omega)$, and B_t is the $M \times N$ -dimensional Brownian motion. Then there exists a constant γ_0 and a set $A \subset \Omega$ such that as $\gamma > \gamma_0$, we have

$$\lim_{t \rightarrow \infty} P(Z_t \in A) = 1,$$

where A is the set of global minima of U .

Proof. Let

$$s = s(t) = \int_0^t \eta(u) du \quad (13)$$

denote its inverse function as $t = t(s)$. Define

$$Y_s = Z_t = Z_{t(s)}.$$

Then the Eq. (12) becomes

$$dY_s = b(Y_s) ds + \frac{\gamma \sqrt{\eta(t)}}{\sqrt{\log s}} dB_t. \quad (14)$$

In terms of the self-similarity property of the Brownian motion and $\eta(t) dt = ds$, we derive that

$$d\tilde{B}_s := \sqrt{\eta(t(s))} dB_{t(s)} \sim N(0, ds \cdot I),$$

where I is the $(M \times N) \times (M \times N)$ unit matrix. Hence \tilde{B}_s is still a standard Brownian motion on $\mathbb{R}^{M \times N}$. Now Eq. (12) becomes

$$dY_s = b(Y_s) ds + \frac{\gamma}{\sqrt{\log s}} d\tilde{B}_s. \quad (15)$$

From the condition of the present theorem, we see that

$$s(t) \rightarrow \infty \quad \text{as } t \rightarrow \infty$$

and

$$t(s) \rightarrow \infty \quad \text{as } s \rightarrow \infty.$$

Therefore, we have

$$\lim_{t \rightarrow \infty} P(Z_t \in F) = \lim_{s \rightarrow \infty} P(Y_s \in F)$$

for any measurable subset F of \mathbb{R}^M .

By theorems of [2], we deduce that there is a positive constant γ_0 such that as $\gamma > \gamma_0$

$$\lim_{t \rightarrow \infty} P(Z_t \in A) = \lim_{s \rightarrow \infty} P(Y_s \in A),$$

where A is the set of the minima of U as $b(x) = -\text{grad } U(x)$. \square

Remark 1. In Theorem 6, γ_0 could be (roughly) chosen to equal to

$$\sqrt{2 \left(\sup_{x \in \Omega} U(x) - \inf_{x \in \Omega} U(x) \right)}.$$

Remark 2. If there is no energy function U for the dynamics, the action functional defined by

$$A(x, y) = \inf_{\phi} \left\{ S_{0,T}(\phi); S_{0,T}(\phi) = \frac{1}{2} \int_0^T |\phi'_t - b(\phi_t)|^2 dt, \phi \in C_{[0,T]}(\mathbb{R}^M), \right. \\ \left. \phi_0 = x, \phi_T = y, \forall T \geq 0 \right\}$$

could be used to replace U and

$$\gamma_0 = \sqrt{2 \sup_{x,y \in \Omega} A(x,y)}.$$

Similar results as in the above theorem are still true, see [2,14].

Remark 3. Our approach also yields a conclusion which is already noted in [24] at p. 259. When $b(x) = -x$ in Eq. (4) it is pointed out in [24], at p. 259, that the second condition (II) of Section 1, i.e.

$$\int_0^\infty \eta^2(u) du < \infty$$

can be replaced by a much weaker condition

$$\lim_{t \rightarrow \infty} \eta(t) = 0$$

and the conditions

$$\int_0^\infty \eta(u) du = \infty, \quad \lim_{t \rightarrow \infty} \eta(t) = 0$$

are necessary and sufficient for X_t to converge to 0. In fact, our approach also rigorously yields this result. Consider the equation of X_t

$$dX_t = \eta(t) [b(X_t) dt + dB_t].$$

After taking the new time scaling s (see the proof of Theorem 6), we yield that

$$dY_s = b(Y_s) ds + \sqrt{\eta(s)} d\tilde{B}_s$$

if $\eta(t) = \eta(t(s)) \rightarrow 0$ and $U(x)$ has only one minimum, say x_0 (the case considered in [24]; $x_0 = 0$), we know that $X_t \rightarrow x_0$, a.s. This proves the sufficiency. The necessary condition is obvious since if $\eta(t)$ does not go to zero, Y_s will certainly not stay at x_0 at all.

Remark 4. We can of course choose a family of VRLPs decreasing more slowly and at the same time to ensure that the conclusions of Theorem 6 are still true. For example, if we set

$$\beta(t) = \frac{\gamma}{\sqrt{\eta(t) \log \log \int_0^t \eta(u) du}}$$

then we still have the conclusions of Theorem 6.²

3.2. A special case

In some situations, it is not possible to separate the drift term b from the Brownian motion B_t . And sometimes the data sent as an input to the network is noise-con-

² We thank T. Heskes and H.J. Kappen for pointing out this remark.

taminated also. This is equivalent to ask if there exists a family of $\eta(t)$ such that X_t converges to the global minima of U , where X_t is the solution of

$$dX_t = \eta(t)(b(X_t) dt + dB_t).$$

From Theorem 6, we know that above requirement is equivalent to say that for $t \geq 1$,

$$\eta(t) \log \int_0^t \eta(u) du = \gamma^2$$

or

$$\log \int_0^t \eta(u) du = \frac{\gamma^2}{\eta(t)}.$$

Differentiating on both sides of the equation above, we have

$$\frac{\eta(t)}{\int_0^t \eta(u) du} = - \frac{\eta'(t) \gamma^2}{\eta(t)^2}$$

or

$$\eta'(t) = - \frac{\eta^3(t)}{\gamma^2 \int_0^t \eta(u) du}. \quad (16)$$

If we are able to solve the above equation and prove that its solution satisfies the conditions of Theorem 6, we obtain a family of new VLRLPs $\eta(t)$. However, it seems that it is not easy, at least theoretically, to find a solution of the above equation and until today we do not know very much about the ODE with delay of form (16), see for example [25] and references therein. But $\eta(t)$ could be numerically computed (Fig. 3) and we have the following estimate. In terms of the nonincreasing property of $\eta(t)$, we have

$$- \frac{\eta^2(t)}{\gamma^2 t} \leq \eta'(t) \leq - \frac{\eta^3(t)}{\eta(1) \gamma^2 t}$$

for $t \geq 1$, which implies that

$$\frac{\gamma^2 \eta(1)}{\eta(1) \log t + \gamma^2} \leq \eta(t) \leq \frac{\gamma \eta(1)}{\sqrt{\eta^2(1) \log t + \gamma^2}}$$

for $\gamma \geq \gamma_0$. This also proves that the condition (11) in Theorem 6 for $\eta(t)$ is fulfilled.

By combining Theorem 6 and all conclusions above now, we come to the main theorem of the present paper. We say that a family of VLRLPs is optimal if it guarantees the learning algorithm to converge to the global minima of the energy function.

Theorem 7. *A family of optimal VLRLPs in the learning algorithm with VLRLPs, i.e. in the following stochastic differential equation*

$$dX_t = \eta(t)(b(X_t) dt + dB_t)$$

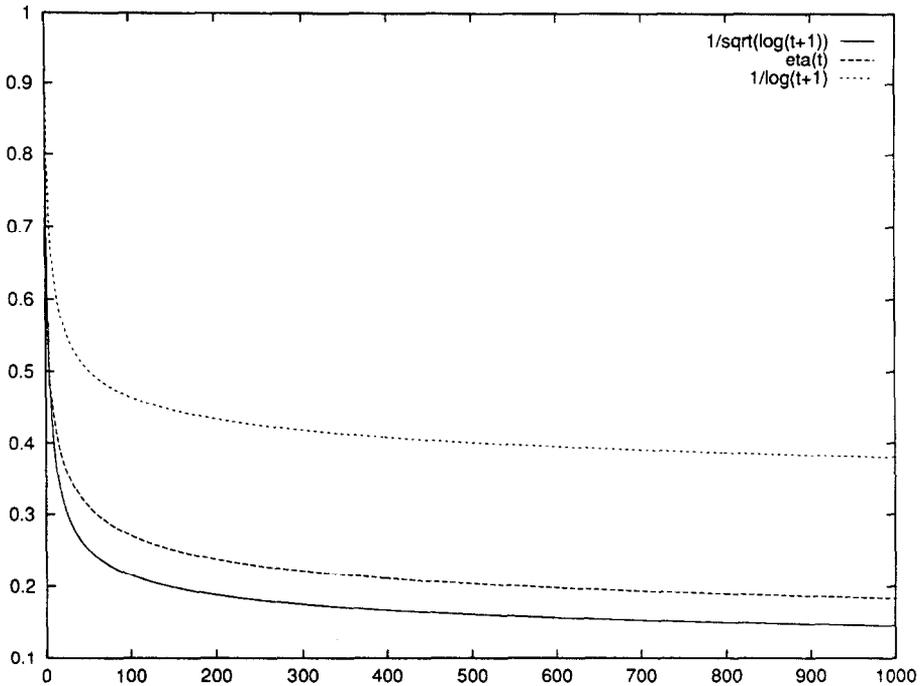


Fig. 3. The function $\eta(t)$ with $\gamma = 1$ defined in Theorem 7 of Section 3.

is the solution of the equation

$$\eta'(t) = - \frac{\eta^3(t)}{\gamma^2 \int_0^t \eta(u) du}$$

with $\eta(1) = \gamma^2 / (\log \int_0^1 \eta(u) du)$ (see Fig. 3). $\eta(t)$ is bounded from below and above:

$$\frac{\gamma^2 \eta(1)}{\eta(1) \log t + \gamma^2} \leq \eta(t) \leq \frac{\gamma \eta(1)}{\sqrt{\eta(1)^2 \log t + \gamma^2}}$$

for some positive constants $\gamma \geq \gamma_0$, where γ_0 is defined as in Theorem 6.

Next we present numerical simulations for a simple model. The reason for us to consider this simple model here is that we can find γ_0 exactly.

Example 5. Let $U(x) = x^4 + x^3 - 4x^2 + x$ (see Fig. 4). We have a numerical comparison of the following three kind of dynamics: the algorithm with the VLRPs of Theorem 7

$$du_t = - \eta(t)(U'(x_t) dt + dB_t); \tag{17}$$

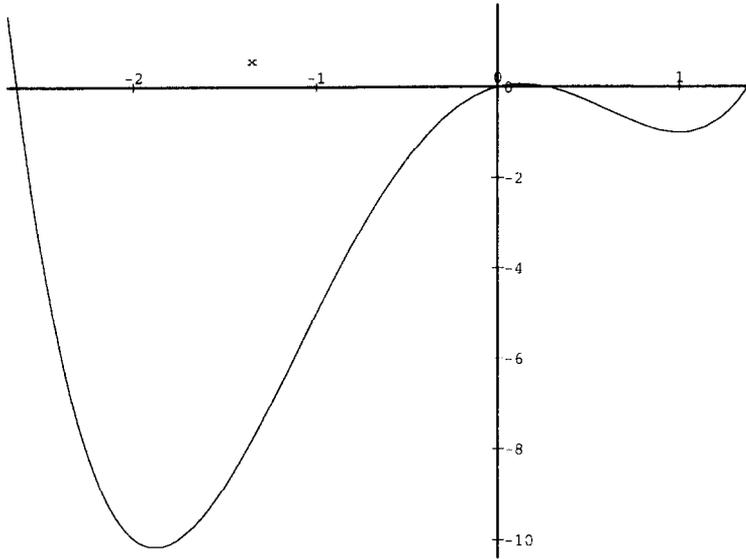


Fig. 4. The potential function $U(x) = x^4 + x^3 - 4x^2 + x$. There are two minima, one is at $x = 1$ (a local minimum) and another is at $x = (-7 - \sqrt{65})/8 = -1.81$ (global minimum).

the algorithm of simulated annealing

$$dv(t) = -U'(v_t) dt + \frac{\gamma}{\sqrt{\log(t+2)}} dB_t; \quad (18)$$

the algorithm with VLRPs of $1/t$

$$dw_t = -\frac{1}{t} (U'(w_t) dt + dB_t). \quad (19)$$

We discretized them with time step $h = 0.01$ (see Eq. (5)) and with initial state $u_0^{(j)} = v_0^{(j)} = w_0^{(j)} = 0.1j - 1$, where $j = 0, \dots, 20$, namely we carried out 21 times simulation with initial state from $[-1, 1]$ for dynamics u_t , v_t and w_t . For each given j after 50000 time iterations we get a solution $u(j)$, $v(j)$ and $w(j)$ corresponding to dynamics (17), (18) and (19), respectively (see Table 1).

Finally, we have

$$u = \frac{\sum_{j=1}^{21} u(j)}{21} = -1.76, \quad v = \frac{\sum_{j=1}^{21} v(j)}{21} = -1.88, \quad w = \frac{\sum_{j=1}^{21} w(j)}{21} = -0.36$$

the average of dynamics (17), (18) and (19) over 21 different initial states. Note that the exact global minima is at $x = -1.81$. The parameter γ is set to 2 (we refer the reader to [2] for an explanation of the choice of this value, γ_0 in Remark 1 is a rather rough choice).

As we expected the dynamics (19) will stay at the local minima with highest probability among dynamics (17), (18) and (19). Dynamics (18) and (19) are more likely to go to the global minima. For the dynamics (18) all 21 simulations are successful in

Table 1

Numerical results of three algorithms (alg.) for initial states from -1 to 1 . Note that only starting from 0.6 the algorithm $u(j)$ fails to arrive to the global minima

Initial state	-1.00	-0.90	-0.80	-0.70	-0.60	-0.50	-0.40	
alg. 1 $u(j)$	-1.87	-1.88	-2.00	-1.95	-1.73	-1.91	-1.83	
alg. 2 $v(j)$	-1.87	-1.79	-2.00	-1.90	-1.83	-1.85	-1.77	
alg. 3 $w(j)$	-1.80	-1.73	-1.78	-1.58	-1.55	-1.58	-1.02	
Initial state	-0.30	-0.20	-0.10	0.00	0.10	0.20	0.30	
alg. 1 $u(j)$	-1.81	-2.14	-1.85	-1.93	-1.91	-1.88	-1.77	
alg. 2 $v(j)$	-1.78	-2.12	-1.87	-1.85	-1.93	-1.88	-1.74	
alg. 3 $w(j)$	-1.27	-1.42	-0.50	-0.02	0.21	0.21	0.25	
Initial state	0.40	0.50	0.60	0.70	0.80	0.90	1.00	Mean value
alg. 1 $u(j)$	-1.90	-1.83	0.97*	-1.97	-1.97	-1.92	-1.82	-1.76
alg. 2 $v(j)$	-1.82	-1.90	-1.92	-1.92	-1.95	-1.87	-1.84	-1.88
alg. 3 $w(j)$	0.68	0.60	0.84	0.97	1.03	0.97	1.02	-0.36

finding the global minima. For dynamics (17) 20 simulations are successful in finding the global minima but one fails. Our numerical results here confirm our theoretical approach.

Finally, a comment should be made about the practical use of the theory presented in this section. Typically, there are two ways to associate the dynamics (4) with a learning algorithm such as the self-organizing Kohonen algorithm, Hebb-type learning, etc. One way is to consider the following learning algorithm

$$\frac{dx_i}{dt} = \eta(t) b(x_i). \tag{20}$$

If we suppose that some stochastic noises are contained in the model, the simplest assumption of it is that now the dynamics (20) takes the form (4). Note that our Theorem 6 and Theorem 7 are proved without any restriction on b , except that $b \in C^1(\Omega)$ (see Remark 2 and [2]), and so it is general enough to cover learning algorithms developed in neural networks.³ In this situation, for avoiding local minima, our approach suggests that it is more reasonable to use the family of vanishing learning rate parameters in Theorem 7 than the one of order $1/t^\alpha$, $1/2 < \alpha \leq 1$. Another way is that the term $\eta(t)B_i$ might be added artificially, following the usual logic of the “annealing” scheme, in order to force the dynamics to jump around until it eventually “settles” near a global minimum. For example for the simple competitive learning defined by Eq. (1) let $b(x) = (E\bar{I}(x, \xi(n+1)))(\xi_j(n+1) - x_{ij})$, $i = 1, \dots, N$, $j = 1, \dots, M$, where E is the expectation with respect to ξ , by adding a noise term $\eta(t)B_i$ to the learning algorithm we assert that the algorithm will reach a global minimum.

³ in [12] it is proved that there is no function U for the self-organizing Kohonen algorithm with the property $b = -\text{grad } U$.

4. Conclusions

Basically, we consider two questions in the present paper. First, a convergence theorem for simple competitive learning algorithm is presented which may be thought of as a replacement of the proof in [17], at p. 222, where the authors claimed that “This result (proof) is, however, somewhat deceptive for two reasons.” Secondly, we rigorously derive a new family of vanishing learning rate parameters for a class of learning algorithms.

Global optimization of learning in neural networks is currently an important subject. How can one be sure that the learning network reaches the optimal state, i.e. the global minimum of some error criterion, and does not get stuck in a local minimum? A well-known strategy to find the global minimum and not just a local minimum is simulated annealing [2], a noise parameter, say temperature, is cooled down slowly. More specifically, we consider the following stochastic differential equation (or Langevin equation):

$$dX_t = -\text{grad } U(X_t) dt + \alpha(t) dB_t, \quad (21)$$

and when

$$\alpha(t) = \frac{\gamma}{\sqrt{\log(t+2)}}, \quad (22)$$

we have

$$\lim_{t \rightarrow \infty} P(X_t \in A) = 1,$$

where A is the set of global minima of U and γ is a constant depending on U .

Learning in neural networks such as the self-organizing Kohonen algorithm, Hebb learning, etc., are also a stochastic process. At each learning step, a training pattern is drawn at random from the environment (the total set of training patterns) and presented to the network. A large learning parameter leads to large fluctuations in the networks representation. So, in a way, the learning parameter can be viewed as a noise parameter akin to the temperature in simulated annealing. A typical case of such learning algorithms (see final chapter of previous section) is

$$dY_t = \eta(t)(b(Y_t) dt + \beta(t) dB_t), \quad (23)$$

a dynamics studied in stochastic approximation theory for many years. Note that when $b = -\text{grad } U$, $\eta(t) = 1$ and $\alpha(t) = \beta(t)$, we have $X_t = Y_t$ and so the case for simulated annealing is just a special case of Eq. (23).

In the present paper we derive a family of vanishing learning rate parameters based upon a rigorous analysis on Eq. (23) and our previous results of simulated annealing in [2]. The new family of vanishing learning rate parameters satisfy the following condition

$$\int_0^\infty \eta(u) du = \infty$$

$$\beta(t) = \frac{\gamma}{\sqrt{\eta(t) \int_0^t \eta(u) du}}, \quad (24)$$

which in general violates the condition (10) found in stochastic approximation theory. Again, we want to point out here that when $\eta(u) = 1$, the rate (22) found in simulated annealing algorithm defined by Eq. (21) is exactly a special case of our results here.

Finally, we like to comment on further possible developments of our results here. Obviously, a case to case and systematic numerical simulations for algorithms developed in neural networks with VLRPs in Theorem 6 and Theorem 7 are quite interesting and is one of our further topics. Theoretically simulated annealing of form (21) has been well studied [2] and, on the other hand, the stochastic approximation theory taking into account the dynamics (23) has developed into a mature field already. In particular, many estimates on convergence rate (in neural networks, convergence rate is called learning error and generalization error) for both algorithms have been established already. We believe that the method developed in this paper serves as a bridge between these two fields and will help us to understand more deeply the behavior of learning algorithms in neural networks and may provide a theoretical basis for the design of practical algorithms that lead to global optimization of learning in neural networks.

Acknowledgements

We would like to thank one of the referees for his extremely careful reading of our manuscript which helped us to considerably improve the presentation.

Appendix A. Proofs of Lemma 3, Theorem 2 and Theorem 5

The tool used in the proof of Theorem 2 and Theorem 5 is the convergence theorem of supermartingale [10,23]. Obviously, here we are not able to provide all elementary facts of martingale theory which is the basis of modern probability theory. But the idea is quite straightforward. We find out the supermartingale (Lyapunov function) related to the simple competitive learning. Based upon the theory of martingale we arrive at all conclusions of Theorem 2 and Theorem 5. We first prove Lemma 3.

A.1. Proof of Lemma 3

“ \Rightarrow ”. First note that if $w_1(0) < w_2(0) < \dots < w_N(0)$, after taking the simple competitive learning, we still have $w_1(n) < w_2(n) < \dots < w_N(n)$, for $n \geq 0$. Suppose that there is a Voronoi tessellation $\Pi \in \Theta$, then there exists

$$z_1 < z_2 < \dots < z_N$$

such that

$$\Pi_i = \left[\frac{z_{i-1} + z_i}{2}, \frac{z_i + z_{i+1}}{2} \right], \quad i = 1, \dots, N,$$

here $(z_0 + z_1)/2 = a$ and $(z_{N+1} + z_N)/2 = b$. So, $y \in \Theta$ implies that $a \leq y_1 < y_2 < \dots < y_N \leq b$.

“ \Leftarrow ”. Trivial. \square

A.2. Proof of Theorem 5

We need to introduce some more notation. Let \mathcal{F}_n be the sigma algebra generated by $\xi(k)$, $k \leq n$, $E(\zeta | \mathcal{F}_n)$ is the conditional expectation for the random variable ζ with respect to the sigma algebra \mathcal{F}_n . Next we are going to find a negative bound for the difference:

$$\sum_{i=1}^N \left[E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n) - |w_i(n) - w_i|^2 \right],$$

and from it we get that $E((w_i(n+1) - w_i)^2)$ is a supermartingale. The supermartingale is the generalization of the Lyapunov function in the deterministic case and the remark after the Example 1 of Section 2.2 is based on this argument.

According to the definition of the algorithm, we have

$$\begin{aligned} & E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n) - |w_i(n) - w_i|^2 \\ &= E(|w_i(n+1)|^2 | \mathcal{F}_n) - 2w_i \cdot E(w_i(n+1) | \mathcal{F}_n) + |w_i|^2 - |w_i(n) - w_i|^2 \\ &= 2\eta(n)(w_i(n) - w_i) \cdot E((\xi(n+1) - w_i(n))I(w_i(n), \xi(n+1)) | \mathcal{F}_n) \\ &\quad + \eta^2(n) E(|\xi(n+1) - w_i(n)|^2 I(w_i(n), \xi(n+1)) | \mathcal{F}_n). \end{aligned}$$

Since $w_i(n)$ and $\xi(n+1)$ are in the set

$$\{|\xi(n+1) - w_i(n)| \leq |\xi(n+1) - w_j(n)|, j \neq i\}$$

if and only if

$$\xi(n+1) \in \Pi(w(n))_i,$$

for $w(n) = (w_i(n), i = 1, \dots, N)$, $w_i(n)$ and $\xi(n+1)$ are independent, we yield that

$$\begin{aligned} & \eta(n)(w_i(n) - w_i) \cdot E((\xi(n+1) - w_i(n))I(w_i(n), \xi(n+1)) | \mathcal{F}_n) \\ &= \eta(n)(w_i(n) - w_i) \cdot \int_{\Pi(w(n))_i} (x - w_i(n))f(x) dx \end{aligned} \tag{25}$$

and

$$\begin{aligned} & \eta(n)^2 E(|\xi(n+1) - w_i(n)|^2 I(w_i(n), \xi(n+1)) | \mathcal{F}_n) \\ &= \eta^2(n) \int_{\Pi(w(n))_i} |x - w_i(n)|^2 f(x) dx. \end{aligned} \tag{26}$$

Furthermore, if we replace the time n in equality (25) and (26) by the stopping time $\sigma_n := \tau(\epsilon) \wedge n = \min(n, \tau(\epsilon))$ all equalities hold. From the definition of the stopping time and the condition (2) in the Theorem 2, we see that

$$\begin{aligned}
 &g(w_1(\sigma_n), \dots, w_N(\sigma_n); w_1, \dots, w_N) \\
 &= \sum_{i=1}^N (w(\sigma_n)_i - w_i) \cdot \int_{\Pi(w(\sigma_n)_i)} (x - w(\sigma_n)_i) f(x) dx \\
 &\leq -h(\epsilon) \\
 &< 0,
 \end{aligned} \tag{27}$$

for a number $h(\epsilon)$ depending only on ϵ . By the condition (3) of Theorem 2, for n large enough, the sign of the term

$$\begin{aligned}
 &\eta(n) \sum_{i=1}^N (w(\sigma_n)_i - w_i) \cdot \int_{\Pi(w(\sigma_n)_i)} (x - w(\sigma_n)_i) f(x) dx \\
 &+ \eta^2(n) \sum_{i=1}^N \int_{\Pi(w(\sigma_n)_i)} |x - w(\sigma_n)_i|^2 f(x) dx
 \end{aligned}$$

is determined by the sign of the following term

$$\begin{aligned}
 &g(w_1(\sigma_n), \dots, w_N(\sigma_n); w_1, \dots, w_N) \\
 &= \sum_{i=1}^N (w(\sigma_n)_i - w_i) \cdot \int_{\Pi(w(\sigma_n)_i)} (x - w(\sigma_n)_i) f(x) dx,
 \end{aligned}$$

and so is negative and we denote it $-h_1(\epsilon) < 0$. This explain the reason why we introduced the function g in Section 2. Without loss of generality, we assume that Eq. (27) is true for $n \geq 1$.

We consider again the term

$$\sum_{i=1}^N |w_i(\sigma_n) - w_i|^2 + \sum_{k=1}^{\sigma_n-1} h_1(\epsilon) \eta(k).$$

After repeating the same argument as before, we conclude that it is still a non-negative supermartingale and so is

$$\sum_{i=1}^N |w_i(\sigma_n) - w_i|^2 + \sum_{k=1}^{\sigma_n-1} h_1(\epsilon) \eta(k).$$

By the convergence of the supermartingale, the limit of

$$\sum_{i=1}^N |w_i(\sigma_n) - w_i|^2 + \sum_{k=1}^{\sigma_n-1} h_1(\epsilon) \eta(k)$$

and

$$\sum_{i=1}^N |w_i(\sigma_n) - w_i|^2$$

are both finite almost surely. Thus,

$$\lim_{n \rightarrow \infty} \sigma_n = \lim_{n \rightarrow \infty} \tau(\epsilon) \wedge n < B$$

almost surely for an integer B satisfying

$$\sum_{k=1}^B \eta(k) h_i(\epsilon) > N \max_{x,y \in \Omega} |x-y|^2 \geq \sum_{i=1}^N |w_i(n) - w_i|^2, \quad \forall n,$$

which implies

$$\tau(\epsilon) < B$$

almost surely. Note that the random time $\tau(\epsilon)$ is bounded by a deterministic quantity B . \square

A.3. Proof of Theorem 2

In terms of the proof of Theorem 5 we see that

$$\begin{aligned} & \sum_{i=1}^N \left[E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n) - |w_i(n) - w_i|^2 \right] \\ & \leq \eta(n) g(w(n), w) + \eta(n)^2 g_1(w(n)) \end{aligned} \quad (28)$$

for

$$g_1(w(n)) = \sum_{i=1}^N \int_{H(w(n), i)} |x - w_i(n)|^2 f(x) dx$$

and $w(n) = (w_1(n), w_2(n), \dots, w_N(n))$. Since $g_1(w(n))$ is uniformly bounded by a constant A the inequality (28) thus becomes

$$\begin{aligned} & \sum_{i=1}^N \left[E(|w_i(n+1) - w_i|^2 | \mathcal{F}_n) - |w_i(n) - w_i|^2 \right] \\ & \leq \eta(n) g(w(n), w) + \eta(n)^2 A \\ & \leq \eta(n) g(w(n), w) + \eta(n)^2 A \left(1 + \sum_{i=1}^N |w_i(n) - w_i|^2 \right). \end{aligned} \quad (29)$$

In terms of the Theorem 7.1 in [23], at p. 43, together with Theorem 5 of the present paper, we arrive at the conclusions of Theorem 1. \square

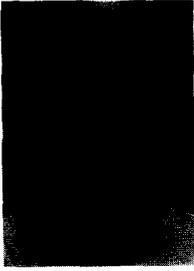
References

- [1] S. Amari, Information geometry of the EM and em algorithms for neural network, METR 94-04, Mathematical Engineering Section, Department of Mathematical Engineering and Information Physics, Faculty of Engineering, The University of Tokyo, Bunkyo-Ku, Tokyo, Japan, 1994.
- [2] S. Alberverio, J.F. Feng and M.P. Qian, Role of noises in neural networks, *Physical Review E* 52 (1995) 6593–6606.

- [3] S. Albeverio, N. Krüger and B. Tirozzi, An extension of Kohonen phonetic maps for speech recognition SFB 237-preprint, Nr. 181, Institut für Mathematik, Ruhr-Universität Bochum, D-44780 Bochum, FRG 1995.
- [4] D. Amit, *Modeling Brain Function* (Cambridge University Press, 1989).
- [5] C. Bishop, *Neural Networks for Pattern Recognition* (Oxford University Press, 1995).
- [6] A. Benveniste, M. Métivier and P. Priouret, *Adaptive Algorithms and Stochastic Approximations* (Springer-Verlag, Berlin, 1990).
- [7] C. Bouton and G. Pagès, Self-organization and convergence of the one-dimensional Kohonen algorithm with non uniformly distributed stimuli, *Stoch. Proc. Appl.* 47 (1993) 249–274.
- [8] C. Bouton and G. Pagès, Convergence in distribution of the one-dimensional Kohonen algorithms when the stimuli are not uniform, *Adv. Appl. Prob.* 26 (1994) 80–103.
- [9] M. Cottrell and J.C. Fort, Etude d'un algorithme d'auto-organisation, *Ann. Inst. H. Poincaré* 23 (1986) 1–20.
- [10] S.Y. Chow and H. Teicher, *Probability Theory* (Springer-Verlag, New York, 1988).
- [11] E. Erwin, K. Obermayer and K. Schulten, Self-organizing maps: Stationary states, metastability and convergence rate, *Biol. Cybern.* 67 (1992) 35–45.
- [12] E. Erwin, K. Obermayer and K. Schulten, Self-organizing maps: Ordering, convergence properties and energy functions, *Biol. Cybern.* 67 (1992) 47–55.
- [13] J.F. Feng, H. Pan and V.P. Roychowdhury, On neurodynamics with limiter function and Linsker's developmental model, *Neural Computation* 8 (1996) 1003–1019.
- [14] J.F. Feng and M.P. Qian, Two-stage annealing in retrieving memories I, in: A. Badrikian, P-A. Meyer and J-A. Yan (eds.), *Probability and Statistics* (World Scientific, Singapore, 1993) 149–176.
- [15] M.I. Freidlin and A.D. Wentzell, *Random Perturbations of Dynamic System* (Springer-Verlag, Berlin, 1984).
- [16] P. Hall and C.C. Heyde, *Martingale Limit Theory and its Application* (Academic Press, New York, 1980).
- [17] J.A. Hertz, A. Krogh and R. Palmer, *Introduction to the Theory of Neural Computation* (Addison-Wesley, Redwood City, California, 1991).
- [18] T. Kohonen, *Self-Organization and Associative Memory (3rd ed.)* (Springer-Verlag, Berlin, 1989).
- [19] H. Kushner and D. Clark, *Stochastic Approximation Methods for constrained and Unconstrained Systems* (Springer-Verlag, New-York, 1978).
- [20] W. Konen, T. Maurer and C. von der Malsburg, A fast dynamic link matching algorithm for invariant pattern recognition, Technical report, Institut Für Neuroinformatik IR-INI 93-05, Ruhr-Universität Bochum, D-44780 Bochum, FRG 1993.
- [21] M. Lades, J.C. Vorbrüggen, J. Buhmann, J. Lange, C. von der Malsburg, R.P. Würtz and W. Konen, Distortion invariant object recognition in the dynamic link architecture, *IEEE Transactions on Computers* 42 (1993) 300–311.
- [22] T. Martinetz and K. Schulten, Topology representing networks, *Neural network* 7 (1994) 507–522.
- [23] M.B. Nevel'son and R.Z. Has'minskii, Stochastic approximation and recursive estimation, Translation of Math. Monograph 47, Amer. Math. Soc., Providence, Rhode Island, 1976).
- [24] H. Ritter, T. Martinetz and K. Schulten, *Neural Computation and Self-Organizing Maps* (Addison-Wesley Publishing Company, Reading, 1992).
- [25] W. Wischert, A. Wunderlin, A. Pelster, M. Olivier and J. Gros Lambert, Delay-induced instabilities in nonlinear feedback systems, *Physical Review E* 49 (1994) 203–219.



J. Feng received his Ph.D. degree in Probability and Statistics in 1991, Peking University. In 1993 he became an Associate Professor for Probability and Statistics, Peking University. After ten months staying at Tübingen University (Germany) as a visiting scholar, he spent a six months visit at Rome University. Now, he works at the Biomathematics Lab. of the Brabham Institute in Cambridge, UK. His current research interests include theoretical foundation of neural networks, Markov processes and random fields.



B. Tirozzi is Professor of Mathematical Physics at the Department of Mathematics, Rome University "La Sapienza". He obtained a Laurea degree in Physics in 1967, after he did research on elementary particle physics. In 1971 he started to work on Statistical Physics, namely on rigorous results concerning Ising and Heisenberg models, dynamics of a classical system with infinitely many particles, methods of probability theory applied to different topics of statistical mechanics. From 1988 he has been involved in the field of neural networks both from the point of view of rigorous results and of the practical problems. He developed a rigorous theory for the spin glass model and the Hopfield model. Many results of his research have been published in the *Journal of Statistical Physics*, *Communication on Mathematical Physics*, *Nuclear Physics B*, *Helvetica Physica Acta*, and *Nuovo Cimento*.