# GENERALIZED HOPFIELD MODELS FOR NEURAL NETWORKS*

By Th. W. Ruijgrok

Institute for Theoretical Physics, Utrecht**

and A. C. C. Coolen

Department of Medical and Physiological Physics

We present an introductory discussion of the Hopfield model for neural networks, together with a new interpretation of Hebb's learning rule. Using a master equation and the language of spin systems, we then derive an equation describing the approach of any initial pattern towards one of the built in patterns. We conclude with some generalizations like spatial structure and the problem of how to perform invariant pattern recognition.

PACS numbers: 12.90.+b

## 1. Introduction

In spite of the fact that the elements of the human brain, the neurons, are quite well understood and also their interconnection and coupling to the outside world are not insurmountable problems, it is a great mystery to understand how a person, governed by this control centre, can perform all tasks of life and even be selfconscious. Being aware of the gigantic size of the problem we must reduce our ambition to a minimum, which is done by studying a network of units, each of which can be in two states, like a neuron which is either firing or in a quiescent state. We will call these units neurons or also spins, although they have very little in common with any of them.

The effect of neuron $j$ on neuron $i$ is described in terms of the synaptic connection $J_{ij}$. It is excitatory if $J_{ij} > 0$ and inhibitory if $J_{ij} < 0$. Its value will be dictated, as we will show later, by what we want the network to do. The system becomes a dynamical system

---

by letting the spins flip with a certain rate, which depends on all other spins and on the amplitude of spontaneous fluctuations, described by a parameter called temperature. The question is whether the system will always evolve towards the same equilibrium state, or whether perhaps the final state depends on the initial spin pattern. In the latter case the network can be used as a pattern recognition device, where the many (meta)stable equilibria correspond to the patterns stored in memory.

In order to answer this question we compare the network with a spin glass, about which much more is known [1, 2]. The alloy AuFe with one percent iron is a known example of a spin glass. The Au and Fe atoms are fixed on the lattice points, but with a random distribution of the Fe sites. The spins of the Fe atoms are the only degrees of freedom. The metal consists of so many macroscopic subsystems, each with its own random distribution of Fe atoms, that a single measurement of a macroscopic quantity amounts in fact already to a series of measurements for a large number of different Fe substitutions. Because of this self averaging a new measurement using another piece of spin glass will not give new results. The microscopic structure of a spin glass is furthermore characterized by the fact that the interaction between the spins of two Fe atoms can be ferro- as well as anti-ferromagnetic, depending on their distance. As a consequence, for given positions of all Fe atoms, it is impossible to orient all spins in such a way that every pair has the lowest possible energy. This so called frustration causes the groundstate to be highly degenerate or leads to the existence of many metastable states.

Since spin glasses and neural networks are described by similar models we expect that it will be possible to obtain many metastable states in networks (necessary to store many patterns) by letting the system be highly frustrated, i.e., by including both excitatory and inhibitory interactions. Also self averaging will have to be an aspect of the model in order to describe the action of the network in terms of macroscopic quantities. This is true despite the fact that the number of a man's neurons is small compared to Avogadro's number: there are as many Fe atoms in one cc of AuFe as there are neurons in the brains of the whole world population.

We will now show how our artificial neurons are constructed and how they are used to build an active network. In describing a single neuron $j$ ($j = 1, 2, ..., N$) we must distinguish its state, given by $s_j = +1$ or $s_j = -1$, from its structure, given by a $p$-dimensional vector $\vec{\xi}_j = (\xi_j^{(1)}, ..., \xi_j^{(p)})$, where $p$ is the number of patterns we eventually want to store in the network. The instantaneous network state is $\vec{s} = (s_1, ..., s_N)$ of which we hope that it will approach one of the desired patterns close enough. We further imagine that every neuron is equipped with $p$ transmitters $t^{(1)}, ..., t^{(p)}$, where $t^{(\mu)}$ can receive and broadcast in a certain $\mu$-band. Every transmitter is in one of two possible states, i.e., either in $\xi_j^{(\mu)} = +1$ or $\xi_j^{(\mu)} = -1$. For each neuron these states are installed once and for all in a learning phase, never to be changed again. The message $m_j^{(\mu)}$ broadcast by $t^{(\mu)}$ depends on the state $\xi_j^{(\mu)}$ of the transmitter and on the state $s_j$ of the neuron. It is equal to $m_j^{(\mu)} = \xi_j^{(\mu)} s_j$ and therefore equal to $+1$ or $-1$. The $\mu$-band transmitters on all other neurons receive this message and turn it into a contribution to the post synaptic potential on their neurons with a magnitude $1/N$. The sign of this contribution is equal to the product $\xi_i^{(\mu)} m_j^{(\mu)}$, so that the total post synaptic potential of neuron $i$ is obtained by adding the contributions from all neurons

in all bands

$$h_i = \frac{1}{N} \sum_{j=1}^{N} \sum_{\mu=1}^{p} \xi_i^{(\mu)} \xi_j^{(\mu)} s_j = \sum_{j} J_{ij} s_j, \tag{1.1}$$

where

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^{p} \xi_i^{(\mu)} \xi_j^{(\mu)}. \tag{1.2}$$

This is exactly the form of the interaction as it was proposed originally by Hopfield [3]. It goes without saying that it would be a ridiculous oversimplification to identify the transmitters of our artificial neurons with the synapses of real neurons. The suggestion is, however, that, if with the form (1.2) of $J_{ij}$ (or a similar non-symmetric expression) some realistic aspects of neuronal systems can be described, it may be worthwhile to look for a physiological counterpart of the numbers $\xi_j^{(\mu)}$ characterising individual neurons.

When for a given network state $\vec{s}$ the action potentials $h_i$ have been evaluated, these will determine the change of $s_i$ in the next time step. If the effects of fluctuations can be neglected ($T = 0$) the rule is that the new $s_i$ will be parallel to the 'magnetic field' $h_i$. In other words: in a Monte Carlo process a proposed spin flip will be accepted if in doing so the total energy

$$E(\vec{s}) = -\tfrac{1}{2} \sum_i h_i s_i = -\tfrac{1}{2} \sum_{i \neq j} J_{ij} s_i s_j \tag{1.3}$$

will be lowered. Using Eq. (1.2) this energy can also be written as

$$E(\vec{s}) = -\tfrac{1}{2} N \sum_{\mu=1}^{p} q_\mu^2 = -\tfrac{1}{2} N |\vec{q}|^2, \tag{1.4}$$

where [4]

$$q_\mu(\vec{s}) = \frac{1}{N} \sum_{i=1}^{N} s_i \xi_i^{(\mu)} \tag{1.5}$$

is the overlap between the actual pattern $\vec{s}$ and the pattern $\vec{\xi}^{(\mu)} = (\xi_1^{(\mu)}, ..., \xi_N^{(\mu)})$ built from the $\mu$-band transmitter states of all neurons. The spin state $\vec{s}$ therefore will evolve in such a way as to make the length of the vector $\vec{q}(\vec{s}) = (q_1(\vec{s}), ..., q_p(\vec{s}))$ as large as possible, under the condition (1.5). Assuming that the vectors $\vec{\xi}^{(\mu)}$ ($\mu = 1, ..., p$) are mutually orthogonal, it is clear that $|\vec{q}(\vec{s})|$ can be made equal to one, by taking $\vec{s}$ equal to any of the $2p$ vectors $\vec{\xi}^{(\mu)}$ or $-\vec{\xi}^{(\mu)}$. Since it follows from Eq. (1.5) that $|\vec{q}(\vec{s})| \leqslant 1$, we see that the ground state is at least $2p$-fold degenerate, or $p$-fold if we do not distinguish between a pattern and its negative.

It is now clear how the connections $J_{ij}$ should be chosen, if we want the network to perform as an associative memory. This means that some initial patterns should approach equilibrium states, which must belong to the class of $p$ preassigned patterns. The recipe is to construct the $J_{ij}$ as in Eq. (1.2), with the $p$ vectors $\vec{\xi}^{(\mu)}$ equal to the patterns we want

to retrieve. This is Hebb's rule [5]. For a realistic associative memory we require an initial state to eventually recognise a pattern only if it starts sufficiently close to one of the built in patterns. Otherwise we want it just to walk around aimlessly in $\vec{s}$-space, or at most settle into some nonsense pattern. For the Hopfield model at zero temperature, as described above, this is just what happens. For $p \leqslant 8$ the total domain of attraction of the $p$ patterns has been calculated [6] and it was found that the fraction of $\vec{s}$-space belonging to this domain is always larger than 20%. A more detailed description of the time evolution, also at finite temperatures, will be given in the next Section.

## 2. Image evolution [7]

Instead of trying to calculate the time evolution of each spin state, we will consider the distribution $P_t(\vec{s})$ which gives the probability to find the system in state $\vec{s}$ at time $t$. The time dependence of this distribution is assumed to follow from a master equation

$$\frac{d}{dt} P_t(\vec{s}) = \sum_j w_j(F_j\vec{s})P_t(F_j\vec{s}) - \sum_j w_j(\vec{s})P_t(\vec{s}), \qquad (2.1)$$

where the first term in the right hand side describes transitions into the state $\vec{s}$ (gain term) and the second (loss) term gives transitions away from $\vec{s}$. The transition rate is taken as

$$w_j(\vec{s}) = \tfrac{1}{2}(1 - \tanh(\beta s_j h_j)) \qquad (2.2)$$

which has the form shown in figure 1 ($\beta = 1/T$). The operator $F_j$ is defined by $F_j\vec{s} = (s_1, ..., -s_j, ..., s_N)$.

It is out of the question — and moreover completely uninteresting — to solve Eq. (2.1) for a distribution function of so many variables ($s_1, ..., s_N$). The problem would become much more tractable if we could derive a master equation for the distribution of the macroscopic overlap variables $q_\mu$ defined in Eq. (1.5), which in vector notation can be written as

$$\vec{q}(\vec{s}) = \frac{1}{N} \sum_{j=1}^{N} s_j\vec{\xi}_j. \qquad (2.3)$$

This is a vector in a space of only $p$ dimensions. For the distribution of these variables, given by

$$P(\vec{q}, t) = \sum_{\vec{s}} P_t(\vec{s})\delta(\vec{q} - \vec{q}(\vec{s})) \qquad (2.4)$$

it is indeed possible also to derive a master equation. In order to do so use is made of the index set [8] $I_{\vec{\eta}}$, which is defined as the collection of all neuron indices $j$ out of $\{1, ..., N\}$, for which the transmitter state $\vec{\xi}_j$ is equal to the $p$-dimensional vector $\vec{\eta} = (\eta_1, ..., \eta_p)$, where each $\eta_\mu$ is equal to $+1$ or $-1$. Since there are only $2^p$ different transmitter states, it is clear that the number of indices in each set $I_{\vec{\eta}}$, to be denoted by $|I_{\vec{\eta}}|$, will be very large and on the order of $N/2^p$. If the $p$ patterns $\vec{\xi}^{(\mu)}$ are chosen at random, this number
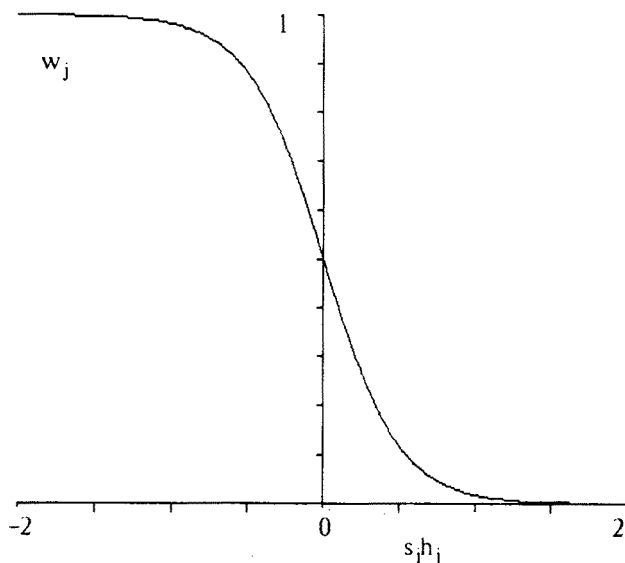
Fig. 1. Transition rate $w_j$ as a function of $s_j h_j$

will be a good estimate for each $|I_{\vec{\eta}}|$. But even if the $\vec{\xi}^{(\mu)}$ are not chosen at random, it will be only in exceptional cases that the numbers $|I_{\vec{\eta}}|$ will deviate appreciably from $N/2^p$. We will therefore ignore the dependence of $|I_{\vec{\eta}}|$ on the patterns $\vec{\xi}^{(\mu)}$ and take $|I_{\vec{\eta}}| = N/2^p$ for all $\vec{\eta}$. With this assumption it is indeed possible to derive the following flow equation

$$\frac{d\vec{q}}{dt} = -\vec{q}(t) + \langle \vec{\eta} \tanh (\beta \vec{\eta} \cdot \vec{q}(t)) \rangle_{\vec{\eta}}. \qquad (2.7)$$

Because of the assumption $|I_{\vec{\eta}}| = N/2^p$ Eq. (2.7) does not depend on the patterns $\vec{\xi}^{(\mu)}$. Some dependence on the patterns occurs when the question is asked how a certain spin state $\vec{s}_0$ at time $t = 0$ evolves. The vector $\vec{s}(t)$ should satisfy the equations

$$q_\mu(t) = \frac{1}{N} \sum_{j=1}^{N} s_j(t) \xi_j^{(\mu)} \quad (\mu = 1, ..., p), \qquad (2.8)$$

where $q_\mu(t)$ is the solution of Eq. (2.7). Since, however, $\vec{s}(t)$ has $N$ components, the $p$ equations (2.8) are absolutely inadequate to determine this $\vec{s}(t)$. The patterns $\vec{\xi}^{(\mu)}$ and the initial spin state $\vec{s}_0$ therefore enter the problem only in that they fix the initial $\vec{q}$-state. At later times information about the spin state $\vec{s}(t)$ is quickly lost.

## 3. Generalizations

It is possible to make the connections $J_{ij}$ nonsymmetric in such a way that the advantages of the form (1.2) are maintained. This can be most easily explained by using the artificial neurons of Section 1 with $p$ transmitters residing on each of them. Transmitter $t^{(\mu)}$ on

neuron $j$ is still sending a message $m_j^{(\mu)} = \xi_j^{(\mu)} s_j$ to all other neurons. Now, however, we allow a receiving $t^{(\nu)}$ transmitter on neuron $i$ to multiply the message by a factor $A_{\nu\mu}$, before it contributes to the action potential. In this way we can understand that $J_{ij}$ becomes

$$J_{ij} = \frac{1}{N} \sum_{\nu,\mu} \xi_i^{(\nu)} A_{\nu\mu} \xi_j^{(\mu)} \tag{3.1}$$

which is nonsymmetric if we choose $A_{\nu\mu}$ to be a nonsymmetric $p \times p$ matrix. The derivation of the previous section can be repeated word by word with the result that the equation for the overlap function now becomes

$$\frac{d\vec{q}}{dt} = -\vec{q}(t) + \langle \vec{\eta} \tanh(\beta\vec{\eta} \cdot A\vec{q}(t)) \rangle_{\vec{\eta}}. \tag{3.2}$$

For the case of two patterns we have solved these equations numerically for three different choises of $A_{\mu\nu}$ and three temperatures. The results are shown in figure 2. The following observations are in order.

a) For a high noise level ($T > 1$) there is no retrieval at all.

b) For $A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and $T < 1$ (the case of Section 2) there is full or imperfect retrieval of the built-in patterns or their negatives.
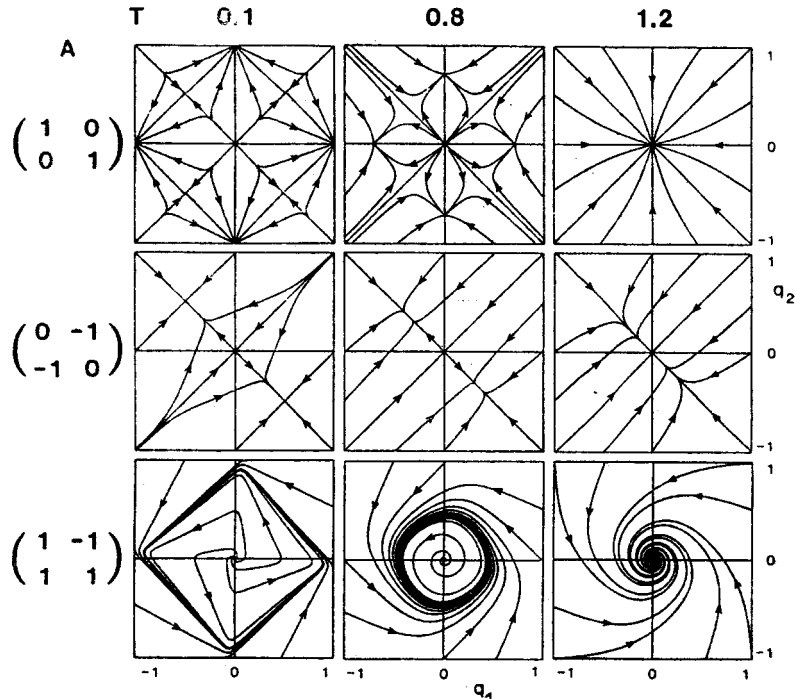


Fig. 2. Examples of the flow described by Eq. (3.2)

c) For $A = \begin{pmatrix} 0 & -1 \\ -1 & 0 \end{pmatrix}$ the vector $\vec{q}$ approaches a mixed pattern.

d) For $A = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$ and $T < 1$ there is a stable limit cycle, which is traversed with long sojourns near the pure patterns.

In another generalization [9] we consider neural networks with arbitrary spatial structure. Let us first assume that the built-in patterns are such that the transmitter states of neighbouring neurons differ only slightly from each other. Then $\vec{\xi}_j$ may be replaced by the $p$-dimensional vector $\vec{\xi}(x)$, where $x$ is the continuously variable position vector of the neuron. As in the previous sections we can again derive an equation for the time dependence of the overlap vector, which now is also space dependent. We find the following differential--integral equation

$$\frac{\partial}{\partial t} \vec{q}(x, t) = \vec{A}(x, \vec{q}) - \vec{q}(x, t), \tag{3.3}$$

with

$$\vec{A}(x, \vec{q}) = \vec{\xi}(x) \tanh (K\vec{\xi}(x) \cdot \int n(x, y)\vec{q}(y, t)dy), \tag{3.4}$$

where $n(x, y)$ represents the density of connections from position $y$ to position $x$. Here $\vec{q}, \vec{A}$ and $\vec{\xi}$ are $p$-dimensional vectors, while $x$ and $y$ may be one-, two- or three-dimensional. $K = \beta J$ is inversely proportional to the temperature ($J$ is simply a constant, introduced to make $K$ dimensionless). If for each $x$ we decompose $\vec{q}(x, t)$ into a component parallel to $\vec{\xi}(x)$ and a component orthogonal to $\vec{\xi}(x)$, so $\vec{q}(x, t) = z(x, t)\vec{\xi}(x) + \vec{q}_\perp(x, t)$, then $\vec{q}_\perp(x, t) = \vec{q}_\perp(x, 0)e^{-t}$. Since this orthogonal component approaches zero and since we are only interested in cases for which the total $\vec{q}(x, t)$ does not vanish after a long time, we can neglect $\vec{q}_\perp(x, t)$. The remaining equation is

$$\frac{\partial}{\partial t} z(x, t) = -z(x, t) + \tanh\left(K \int \hat{n}(x, y)z(x, y)dy\right), \tag{3.5}$$

with

$$\hat{n}(x, y) = n(x, y) (\vec{\xi}(x) \cdot \vec{\xi}(y)). \tag{3.6}$$

Eq. (3.5) is the same as we would have obtained from Eq. (3.3) for only one pattern, except for the difference expressed by Eq. (3.6). In the case that $n(x, y)$ — and therefore also $\hat{n}(x, y)$ — is symmetric, the solution of Eq. (3.5) will always approach some stationary state, no matter what the function $z(x, 0)$ is or how we choose $K$. In order to prove this a free energy functional is defined as

$$F[z] = -\tfrac{1}{2} \int \hat{n}(x, y)z(x, t)z(y, t)dxdy + \frac{1}{K} \int S(z(x, t))dx, \tag{3.7}$$

where the entropy density is

$$S(z) = \tfrac{1}{2}(1+z) \log\left(\frac{1+z}{2}\right) + \tfrac{1}{2}(1-z) \log\left(\frac{1-z}{2}\right). \tag{3.8}$$

We will need the derivative

$$\frac{d}{dz}S(z) = \tfrac{1}{2}\log\frac{1+z}{1-z} \qquad (3.9)$$

and the inverse

$$z = \tanh\left(\frac{d}{dz}S(z)\right). \qquad (3.10)$$

With $T(x, t) \equiv K\int \hat{n}(x, y)z(x, t)dy$, we find for the rate of change of the free energy

$$\frac{dF}{dt} = 1/K\int dx\left(\left(\frac{dS(z)}{dz}\right)_{x,t} - T(x, t)\right)\frac{\partial z(x, t)}{\partial t}.$$

Substitution of Eq. (3.5) and using Eq. (3.10) gives

$$\frac{dF}{dt} = 1/K\int dx\left(\left(\frac{dS(z)}{dz}\right)_{x,t} - T(x, t)\right)\left(\tanh\left(\left(\frac{dS(z)}{dz}\right)_{x,t}\right) - \tan\ (T(x, t))\right). \qquad (3.11)$$

From this it is obvious that $\dfrac{dF}{dt} \leqslant 0$ and that the equality sign holds if and only if $z(x)$ is a stationary solution of Eq. (3.5). For high temperatures, i.e., for small $K$, it is expected that $z(x, t)$ will approach the trivial solution $z = 0$. For lower temperatures, however, it is interesting to investigate whether non zero and possibly non uniform solutions exist.

For that purpose we made a numerical study of the case where $\hat{n}(x, y) = 1$ if $y$ lies in a one- or two-dimensional square of size one with $x$ in its centre, and $\hat{n}(x, y) = 0$ otherwise. In the one-dimensional case we took as boundary conditions $z(-x, t) = z(x, t)$, so that $z(0, t) = 0$, and $z(\infty, t) = \bar{z}$, where $\bar{z}$ is the positive solution (exists if $K > 1$) of $\bar{z} = \tanh(K\bar{z})$. In the two-dimensional case the equation (3.5) was discretized on a $50 \times 50$ lattice with periodic boundary conditions. In both cases the initial state was chosen randomly. The results are shown in figures 3 and 4 for the one-dimensional system and in
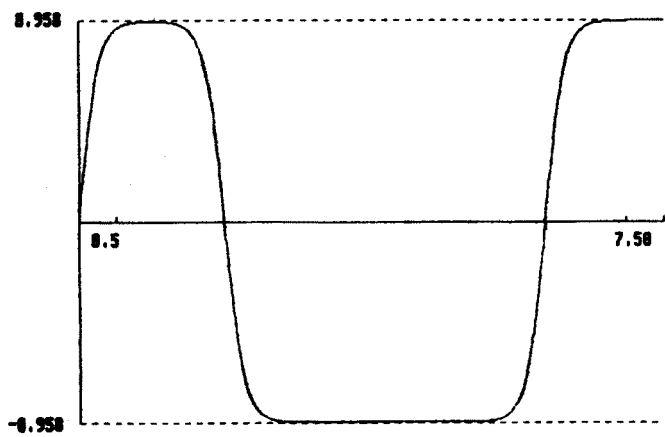


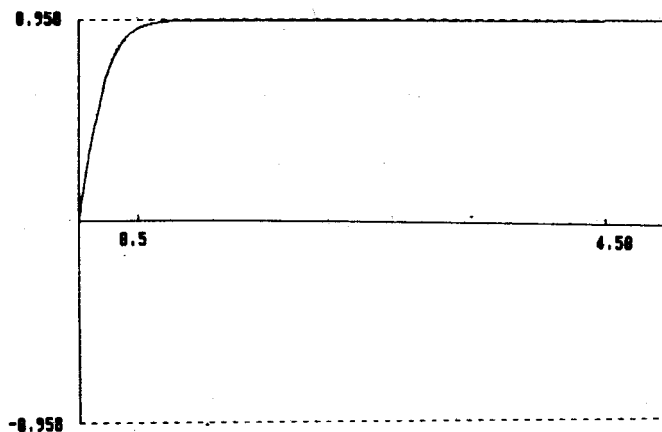Fig. 3. Metastable solution of Eq. (3.5) in one dimension, obtained with $K = 2$

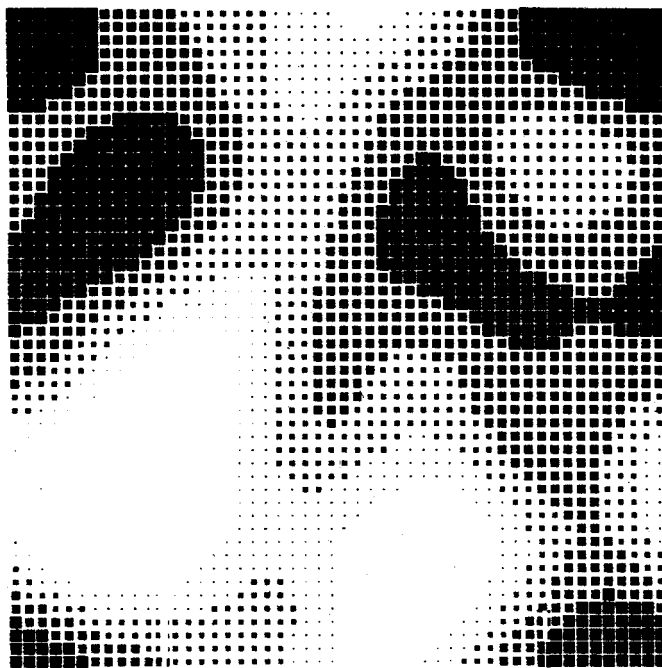Fig. 4. Stable solution of Eq. (3.5) in one dimension, obtained with $K = 2$



Fig. 5. Metastable solution of Eq. (3.5) in two dimensions, obtained with $K = 1.08$

Fig. 5 for the three-dimensional case. In Fig. 3 we see how in a rather short time a metastable state is formed, where each point of the line is either in the state $\bar{z}$ or in the state $-\bar{z}$, with transition regions which are sharper the lower the temperature is. It takes much longer for the metastable states to disappear. The final stable state (for $x > 0$) is shown in Fig. 4. In the two-dimensional case the state shown in Fig. 5 was again reached rather quickly. The accuracy of our calculations was, however, not good enough to see a further

decrease of the free energy, so that we cannot definitely decide whether the state is stable or metastable.

If, on the other hand, all pattern bits $\xi_j^{(\mu)}$ are drawn at random we find Eq. (3.4) being replaced by

$$\vec{A}(x, \vec{q}) = \langle \vec{\eta} \tanh (K\vec{\eta} \cdot \int n(x, y)\vec{q}(y, t)dy)\rangle_{\vec{\eta}}. \tag{3.12}$$

If $\vec{q}$ has only one non-zero component its evolution in time is again given by (3.5), after replacing $\hat{n}(x, y)$ by $n(x, y)$. If $n(x, y)$ is a compact integral operator one can prove [10] that for the system (3.3/3.12) there is a critical temperature $K_c^{-1} \leqslant Mp$, where $M$ is the norm of $n(x, y)$. For translation invariant systems, i.e. $n(x, y) = n(x-y)$, $K_c = 1$. In these calculations $n(x, y)$ is assumed to be normalized such that $\int dy\, n(x, y) = 1$. In general it appears that, for non-zero temperatures, there is a minimum correlation length in the solutions of (3.3/3.12), as they bifurcate from the trivial solution $\vec{q}(x, t) = \vec{0}$, which can be much larger than the range of the connections. This opens up the possibility to store and retrieve complete patterns from local clues, using restricted range connections only. For an interaction function $\hat{n}(x-y)$ with a gaussian shape in three dimensions, for example, the ratio between the minimum correlation length $\Lambda$ and the width $\sigma$ of the interaction kernel is given by:

$$\Lambda(T)/\sigma = \pi \sqrt{2/3} \,(\log (1/T))^{-1/2}. \tag{3.13}$$

Finally we want to address the problem of invariant pattern recognition [11]. One can combine the storage of patterns in a neural network with the storage of a transformation $T$. The latter can be done by choosing connections as follows:

$$J_{ij} = \langle (T\vec{s})_i s_j\rangle. \tag{3.14}$$

If, in principle, $T$ can be written in the form $(T\vec{s})_i = \text{sgn} (\Sigma G_{ij}s_j)$ for some matrix $G$, than (3.14) will in general yield connections $J_{ij}$ such that

$$\text{sgn} (\sum_j J_{ij}s_j) = (T\vec{s})_i. \tag{3.15}$$

The easiest examples are index mappings: $(T\vec{s})_i \equiv s_{\pi(i)}$ for some $\pi: \{1, ..., N\} \to \{1, ..., N\}$. If the spins represent pixels on a screen, then all transformations of the screen onto itself can be written in this way. For these transformations Eq. (3.14) shows that $J_{ij} = \delta_{\pi(i),j}$ so Eq. (3.15) clearly holds. A neural network with zero noise level and synchronous updating of the spins, equipped with the connections (3.14), would show an evolution in time, equivalent to repeated iteration of the transformation $T$ on its microscopic state. If combined with the standard storage of patterns this network will perform pattern recognition, invariant under the transformation group generated by $T$, provided the relative weight of the two contributions to the connection matrix is properly chosen. For the index mappings the optimal choice would be

$$J_{ij} = \frac{1}{N} \sum_\mu \xi_i^{(\mu)}\xi_j^{(\mu)} + \tfrac{1}{2}\delta_{\pi(i),j}. \tag{3.16}$$

## REFERENCES

[1] D. Chowdhury, *Spin Glasses and other Frustrated Systems*, World Scientific, Singapore 1986.

[2] K. Binder, A. P. Young, *Rev. Mod. Phys.* **58**, 801 (1986).

[3] J. J. Hopfield, *Proc. Natl. Acad. Sci. U.S.A.* **79**, 2554 (1982).

[4] D. J. Amit, H. Gutfreund, H. Sompolinsky, *Phys. Rev.* **A32**, 1007 (1985).

[5] D. O. Hebb, *The Organization of Behaviour*, Wiley, New York 1949.

[6] A. C. C. Coolen, H. J. J. Jonker, Th. W. Ruijgrok, to be published.

[7] A. C. C. Coolen, Th. W. Ruijgrok, *Phys. Rev.* **A38**, 4253 (1988).

[8] J. L. van Hemmen, D. Grensing, A. Huber, R. Kühn, *Z. Phys.* **B65**, 53 (1986).

[9] A. C. C. Coolen, J. J. Denier van der Gon, Th. W. Ruijgrok, to be published.

[10] A. C. C. Coolen, J. J. Denier van der Gon, to be published.

[11] A. C. C. Coolen, F. W. Kuijk, to be published.