

# Dynamics of diluted and asymmetric neural network models

Autor(en): **Derrida, B.**

Objektyp: **Article**

Zeitschrift: **Helvetica Physica Acta**

Band (Jahr): **62 (1989)**

Heft 5

PDF erstellt am: **18.04.2024**

Persistenter Link: <https://doi.org/10.5169/seals-116046>

## Nutzungsbedingungen

Die ETH-Bibliothek ist Anbieterin der digitalisierten Zeitschriften. Sie besitzt keine Urheberrechte an den Inhalten der Zeitschriften. Die Rechte liegen in der Regel bei den Herausgebern.

Die auf der Plattform e-periodica veröffentlichten Dokumente stehen für nicht-kommerzielle Zwecke in Lehre und Forschung sowie für die private Nutzung frei zur Verfügung. Einzelne Dateien oder Ausdrucke aus diesem Angebot können zusammen mit diesen Nutzungsbedingungen und den korrekten Herkunftsbezeichnungen weitergegeben werden.

Das Veröffentlichen von Bildern in Print- und Online-Publikationen ist nur mit vorheriger Genehmigung der Rechteinhaber erlaubt. Die systematische Speicherung von Teilen des elektronischen Angebots auf anderen Servern bedarf ebenfalls des schriftlichen Einverständnisses der Rechteinhaber.

## Haftungsausschluss

Alle Angaben erfolgen ohne Gewähr für Vollständigkeit oder Richtigkeit. Es wird keine Haftung übernommen für Schäden durch die Verwendung von Informationen aus diesem Online-Angebot oder durch das Fehlen von Informationen. Dies gilt auch für Inhalte Dritter, die über dieses Angebot zugänglich sind.

DYNAMICS OF DILUTED AND ASYMMETRIC NEURAL NETWORK MODELS

B. Derrida

*Service de Physique Théorique*

*de Saclay*

*Laboratoire de l'Institut de Recherche Fondamentale*

*du Commissariat à l'Energie Atomique*

*F-91191 Gif-sur-Yvette Cedex*

Abstract

This talk is a short review of the dynamical properties of diluted asymmetric neural networks. The time evolution of the projection of a configuration on a stored pattern can be obtained exactly. One can also calculate the distribution of the activities of the neurons and show that the dynamics are chaotic even in the good retrieval phase.

12th Gwatt Workshop, Complex Systems, Oct. 13-15, 1988

## 1. INTRODUCTION

The theory of neural networks has become, in the last few years, a very active field of statistical mechanics (for recent reviews see references<sup>1-9</sup>). One of the reasons why physicists (from statistical mechanics) are interested by the brain is rather obvious. The brain is composed by a large number  $N$  ( $\approx 10^{12}$ ) of neurons which interact through synapses ( $10^2$ - $10^4$  per neuron). It is therefore tempting to try to describe the properties of such a system by the technics of statistical mechanics. Another reason which makes neural networks attractive to physicists is that their time evolution is usually not governed by any hamiltonian. Thus to describe their dynamics, one needs to develop new theoretical methods which should be useful to study all kinds of systems far from equilibrium.

The simplest neural network models which have been considered consist in assuming that the state of each neuron  $i$  at time  $t$  is represented by an Ising variable  $S_i(t)$

$$\begin{aligned} S_i(t) &= +1 \quad \text{if the neuron } i \text{ is firing} \\ S_i(t) &= -1 \quad \text{if the neuron } i \text{ is quiescent} \end{aligned} \tag{1}$$

and that the synapsis  $J_{ij}$  between neuron  $j$  and neuron  $i$  is a real number ( $J_{ij} > 0$  if the synapsis is excitatory and  $J_{ij} < 0$  if it is inhibitory). In general the matrix  $J_{ij}$  is nonsymmetric ( $J_{ij} \neq J_{ji}$ ) because the synapses are non symmetric. One then has to choose a dynamical rule to make the system evolve in time. A simple way consists in saying that at time  $t$  the neuron  $i$  receives a potential  $V_i(t)$  given by

$$V_i(t) = \sum_j J_{ij} S_j(t) \tag{2}$$

and that the state of neuron  $i$  at time  $t + 1$  depends on  $V_i(t)$  in a probabilistic way

$$\begin{aligned} S_i(t+1) &= +1 \quad \text{with probability } f(V_i(t)) \\ S_i(t+1) &= -1 \quad \text{with probability } 1-f(V_i(t)) \end{aligned} \tag{3}$$

$f(x)$  is an increasing function such that  $f(x) \rightarrow 0$  if  $x \rightarrow -\infty$  and  $f(x) \rightarrow 1$  if  $x \rightarrow +\infty$ . For example

$$f(x) = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{x}{T}\right) \quad (4)$$

where  $T$  plays the role of a temperature.

Dynamics like (2-4) are commonly used in Monte Carlo simulations. The main difference with more usual spin models of Statistical Mechanics is that here the matrix  $J_{ij}$  is non symmetric. Therefore there is no hamiltonian, no partition function.

The first property of such neural network models is that they can memorize patterns by choosing properly the synapses  $J_{ij}$ . Assume that we have  $N$  neurons  $S_i = \pm 1$  and we want to store  $p$  patterns  $\{\xi_i^\mu\}$  of  $N$  bits each

$$\begin{aligned} 1^{st} \text{ pattern } \quad \xi_i^{(1)} &= +1 \text{ or } -1 \quad 1 \leq i \leq N \\ &\cdot \\ &\cdot \\ &\cdot \\ p^{th} \text{ pattern } \quad \xi_i^{(p)} &= +1 \text{ or } -1 \quad 1 \leq i \leq N \end{aligned} \quad (5)$$

We will say that pattern  $\{\xi_i^{(\mu)}\}$  is memorized if for the dynamics (2-4) there is an attractor near this pattern. So the problem is to choose the  $J_{ij}$  in order to make the attractors as close as possible to the stored patterns. A simple way of measuring the distance between a spin configuration  $\{S_i(t)\}$  and a pattern is to calculate their overlap

$$m_\mu(t) = \frac{1}{N} \sum_{i=1}^N \xi_i^{(\mu)} S_i(t) \quad (6)$$

There exist several choices of the  $J_{ij}$  which give attractors in the neighbourhood of the stored patterns  $\xi_i^{(\mu)}$ . Some of these choices lead to interesting effects like short on long term memory, forgetting<sup>10-11</sup>. The discussion will be limited here to one of the simplest rules (the Hebb rule<sup>4</sup>) which give an expression of the  $J_{ij}$  in terms of the patterns

$$J_{ij} = \frac{1}{C} \sum_{\mu=1}^p \xi_i^{(\mu)} \xi_j^{(\mu)} \quad (7)$$

where  $C$  is the number of synapses of each neuron (for simplicity one can assume that  $C$  does not depend on  $i$ ).

## 2. THE HOPFIELD MODEL<sup>12,13,4</sup>

As long as the matrix  $J_{ij}$  is non symmetric, it is not easy to use the methods of Statistical Mechanics (Partition function etc ...). The idea of Hopfield was to consider a simpler situation where

$$C = N - 1 \quad (8)$$

i.e. each neuron interacts with each other neuron and the  $J_{ij}$  are given by (7). Then one knows that with the dynamics (2-4) the system will evolve to an equilibrium described by an Hamiltonian  $\mathcal{H}$  at temperature  $T$

$$\mathcal{H}(\{S_i\}) = - \sum_{ij} J_{ij} S_i S_j \quad (9)$$

i.e. each configuration  $\{S_i\}$  is visited in the long time limit with a probability  $P_{eq}(\{S_i\}) = \exp [- \mathcal{H}(\{S_i\})/T]$ .

When one considers the  $J_{ij}$  given by the Hebb rule (7), the  $J_{ij}$  take both positive and negative values and phase space is composed of many valleys like in spin glass problems<sup>14,15</sup>.

Amit, Gutfreund and Sompolinsky<sup>4,16</sup> have studied the equilibrium properties (the thermal equilibrium) in great detail using replica technics which had been developed previously in the study of spin glasses. They found a phase diagram with the following shape :

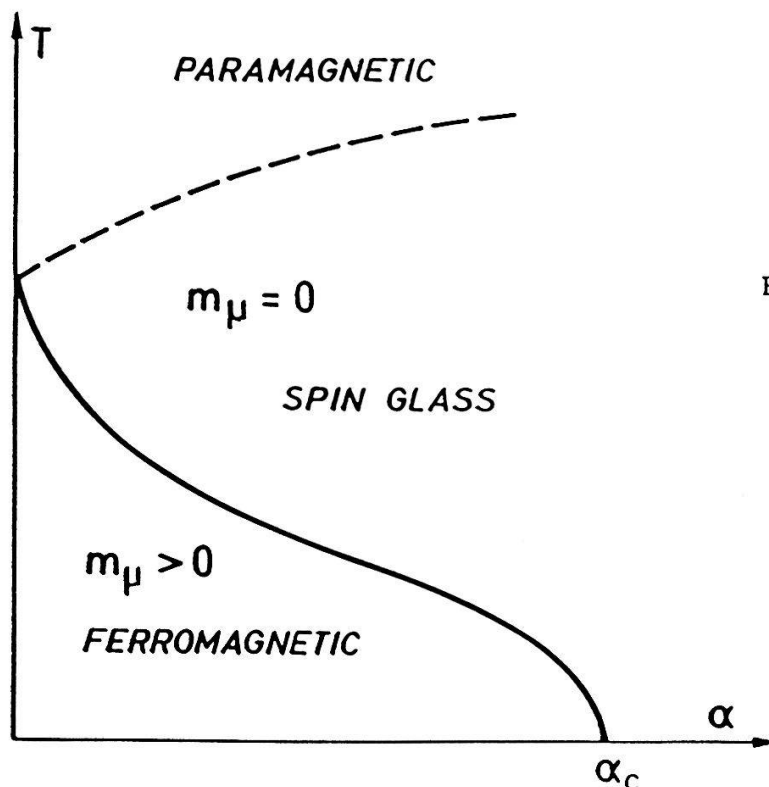


Figure 1.

when the  $p$  patterns are chosen at random (i.e. :  $\xi_i^{(\mu)} = +1$  or  $-1$  with equal probability). The parameter  $\alpha$  which appears in figure 1 is defined as the ratio of the number of stored patterns  $p$  divided by the number of synapses  $C$  per neuron (in the case of the Hopfield model  $C = N - 1$  but this will not be the case for the diluted model discussed in section 3)

$$\alpha = p/C \quad (10)$$

The important line is the phase boundary between the ferromagnetic phase and the spin glass phase. In the ferromagnetic phase, there exists a minimum in the free energy landscape with  $m_\mu > 0$ , i.e. one expects a valley near each pattern  $\xi^{(\mu)}$ . In the spin glass and the paramagnetic phases  $m_\mu = 0$  and therefore this local minimum disappears. The phase diagram obtained by Amit et al<sup>4,16</sup> has more structure than what is shown on figure 1 (transition line where the symmetry of replica is broken, spin glass phase, etc ...) but this will not be discussed here.

We see that with the Hebb rule (7) the system is able to memorize the patterns as long as the number of stored patterns  $p = C \alpha$  does not exceed a certain value  $\alpha_c(T)$ . For  $\alpha > \alpha_c(T)$ , there is a complete deterioration and no pattern is memorized. It turns out<sup>4,16</sup> that the transition from the ferromagnetic phase to the spin glass phase is a first order transition and that  $m_\mu$  has a jump. At  $T = 0$ , one finds that

$$\alpha_c \simeq .14$$

and  $m_\mu$  jumps from a value  $\simeq .95$  to 0. Since the fraction of wrong bits is given by  $\frac{1-m_\mu}{2}$ , we see that up to  $\alpha_c$  the patterns are memorized with very few mistakes.

The calculations done on the Hopfield model can be generalized to various situations (see the reviews<sup>4-5</sup>) by modifying or by replacing the Hebb rule (7) by other rules<sup>10,11</sup>.

There are however several difficulties in the Hopfield model

- (1) The calculations are done at equilibrium (using replica) but one does not know how to describe analytically dynamics.
- (2) The symmetry of the synapses ( $J_{ij} = J_{ji}$ ) is essential in this approach although the synapses are known to be non symmetric in the brain.
- (3) All the neurons are connected ( $C = N - 1$ ) and that too is not realistic.

(4) The forgetting catastrophe : if  $\alpha > \alpha_c$ , i.e. the number of input patterns exceeds a maximal value, all the patterns are forgotten at once.

(5) Various difficulties when one extends the calculations to the case of correlated patterns.

(6) With symmetric interactions, one can expect minima in the free energy landscape but there is no way to memorize temporal sequences of patterns.

### 3. NON SYMMETRIC - DILUTED NETWORKS

It turns out that one can construct a neural network model with non symmetric synapses for which the dynamics can be solved exactly<sup>17</sup>. The model consists of  $N$  neurons  $S_i = \pm 1$  and the synapses  $J_{ij}$  are given by

$$J_{ij} = C_{ij} \sum_{\mu=1}^p \xi_i^{\mu} \xi_j^{\mu} \quad (11)$$

where the  $\{\xi_i^{\mu}\}$  is the  $\mu^{\text{th}}$  pattern and  $C_{ij}$  is a random number which represents the dilution

$$\begin{aligned} C_{ij} &= 1 \quad \text{with probability } C/N \\ C_{ij} &= 0 \quad \text{with probability } 1 - C/N \end{aligned} \quad (12)$$

The  $C_{ij}$  and  $C_{ji}$  are independent random variables and therefore the matrix  $J_{ij}$  is no longer symmetric. The spin still evolve according to the following rules

$$\begin{aligned} S_i(t+1) &= +1 \quad \text{with prob } p_i(t) \\ S_i(t+1) &= -1 \quad \text{with prob } 1 - p_i(t) \end{aligned} \quad (13)$$

where

$$p_i(t) = \frac{1}{2} + \frac{1}{2} \tanh \left( \sum_j J_{ij} S_j(t) / T \right) \quad (14)$$

which gives in the low temperature limit

$$S_i(t+1) = \text{sgn} \left[ \sum_j J_{ij} S_j(t) \right] \quad (15)$$

The situation for which this model can be solved is the limit  $N \rightarrow \infty$ ,  $C$  being finite or infinite with the constraint that

$$C \ll \log N \quad (16)$$

The reason why this condition makes the model soluble would be too long to explain here in detail<sup>15,17</sup>. Let me just say that because the system is very diluted (see eq. 12), the structure of the network is locally a tree. At zero temperature, and for large  $C$ , one gets a simple expression for the time evolution of  $m_\mu(t)$

$$m_\mu(t+1) = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{+\infty} dy e^{-y^2} \text{sign}[m_\mu(t) - y\sqrt{2\alpha}] \quad (17)$$

where  $\alpha = p/C$  ( $p$  is the number of stored patterns). The dynamics are fully described by the map (17). One sees from (17) that there is a critical value  $\alpha_c$  of  $\alpha$

$$\alpha_c = 2/\pi \quad (18)$$

If  $\alpha > \alpha_c$ , the number of patterns memorized is too large and the only attractive fixed point of (17) is  $m_\mu^* = 0$ . The system does not remember anything.

If  $\alpha < \alpha_c$ , there appears an attractive fixed point  $m_\mu^* \neq 0$  of (17) corresponding to the attractor near the pattern  $\mu$ . One should notice that the retrieval is not perfect since  $m_\mu^* \neq 1$  (the fraction of wrong bits is given by  $(1 - m_\mu^*)/2$ ).

In the above calculation, the typical projection of one pattern  $\mu$  on another pattern  $\nu$  was  $N^{-1/2}$  :

$$\frac{1}{N} \sum_i \xi_i^{(\mu)} \xi_i^{(\nu)} \sim N^{-1/2} \quad (19)$$

for all pairs  $\mu$  and  $\nu$ . One can generalize it to describe other situations. If one considers that  $p$  patterns are random but that two of them (patterns 1 and 2) are correlated

$$\frac{1}{N} \sum_i \xi_i^{(1)} \xi_i^{(2)} = Q \quad (20)$$

one can write<sup>17</sup> equations similar to (17) to describe the time evolution of  $m_1(t)$  and  $m_2(t)$ . Because of (20), the time evolution of  $m_1(t)$  and  $m_2(t)$  are



coupled and one finds that there are two critical values of  $\alpha$  :

$$\begin{aligned}\alpha_1 &= \frac{2}{\pi} (1-Q)^2 \\ \alpha_2 &= \frac{2}{\pi} (1+Q)^2\end{aligned}\tag{21}$$

For  $\alpha > \alpha_2$ , the only fixed point is  $m_1^* = m_2^* = 0$ . Too many patterns have been stored. The system does not remember anything.

For  $\alpha_1 < \alpha < \alpha_2$ , there is an attractive fixed point  $m_1^* = m_2^* \neq 0$ . The system remembers patterns 1 and 2 but cannot distinguish them.

For  $\alpha < \alpha_1$ , there is an attractive fixed point  $m_1^* > m_2^*$ . The system can distinguish the two patterns.

There are some limiting cases which can be easily understood.

If  $Q \rightarrow 0$ , the patterns become uncorrelated and  $\alpha_1$  and  $\alpha_2 \rightarrow \alpha_c$ .

If  $Q \rightarrow 1$ ,  $\alpha_1 \rightarrow 0$ . If the 2 patterns become identical, it is impossible to distinguish them.

#### 4. DISTRIBUTION OF ACTIVITIES AND CHAOS IN DILUTED NETWORKS

For systems like the Hopfield model (section II), where the dynamics are governed by an hamiltonian, one knows that at zero temperature, the spin configuration always converges to a fixed point in phase space. The system gets trapped in a local minimum of energy. This means that after a transient part, all the spins  $S_i(t)$  remain fixed in time for ever.

For systems with non symmetric interactions, there is no hamiltonian and the attractors in phase space can be cycles with arbitrarily long periods (at least if the system size is large enough). So even when the system has reached its attractor, the spins  $S_i(t)$  keep changing with time.

For the diluted network defined in section III, it is possible to describe quantitatively the motion of the spins. For simplicity, let us discuss the case of a configuration  $\{S_i(t)\}$  which has a non zero projection  $= m_\mu^*$  (the fixed point of (17)) on a single pattern  $\mu$ .

If one defines the activity  $a_i$  of spin  $i$  as

$$a_i = \lim_{t_0 \rightarrow \infty} \frac{1}{t_0} \sum_{t=1}^{t_0} S_i(t)\tag{22}$$

it is possible<sup>18</sup> to obtain a closed expression for the whole histogram  $P(a)$  of the  $a_i$  :

$$P(a) = \frac{1}{N} \sum_{i=1}^N \delta(a - a_i) \quad (23)$$

The final result can be parametrized in the following way

$$P(a) = \frac{1}{2} \left( \frac{1-q}{q} \right)^{1/2} \exp \left( z^2 - \frac{(z\sqrt{1-q} - m_\mu^*/\sqrt{2\alpha})^2}{q} \right) \quad (24)$$

where

$$a = \text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-x^2} dx \quad (25)$$

$m_\mu^*$  is the fixed point of (17) and  $q$  is the attractive fixed point of

$$q(t+1) = -1 + \frac{2}{\sqrt{\pi}} \int_{-\infty}^{+\infty} e^{-y^2} dy \text{erf} \left( \left| \frac{y\sqrt{1+q(t)} + m_\mu^*/\sqrt{\alpha}}{\sqrt{1-q(t)}} \right| \right) \quad (26)$$

The shape of  $P(a)$  depends on  $\alpha$ . Two typical shapes (both for  $\alpha < \alpha_c$ ) are shown in Fig.2

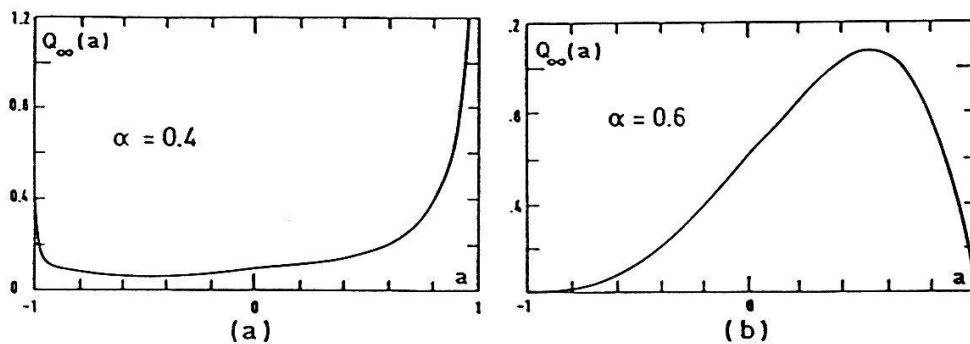


Figure 2.

Fig. 2: The distribution of activities  $P(a)$  for the diluted asymmetric model.

The fact that the distribution  $P(a)$  is a continuous function without

delta peak means that no spin is fixed. They keep moving for ever. This is a major difference between models with only symmetric synapses (for which all the spins are fixed at zero temperature) and models with asymmetric synapses.

Even when the system converges to a single pattern attractor, it is possible to show that the dynamics remain chaotic: if one considers two configurations  $\{S_i(t)\}$  and  $\{\tilde{S}_i(t)\}$  which have both the projection  $m_\mu^*$  on the same pattern  $\mu$  and zero projection on all the other patterns, it is possible to calculate the time evolution of their overlap  $q(t)$ :

$$q(t) = \frac{1}{N} \sum_{i=1}^N S_i(t) \tilde{S}_i(t) \quad (27)$$

It turns out<sup>17</sup> that  $q(t)$  is given by the recursion (26). For  $\alpha < \alpha_c$ , (26) has two fixed points (with  $a > 0$ ):  $q=1$  and  $q^* \neq 1$ . The fixed point  $q=1$  corresponds to two identical configurations  $\{S_i(t)\} = \{\tilde{S}_i(t)\}$ . Clearly if two configurations are identical, they remain identical at later times. This fixed point  $q=1$  is however unstable whereas  $q^*$  is an attractive fixed point. This implies that if two configurations are initially slightly different, their trajectories in phase space diverge and their overlap always converges to  $q^*$ . The fact that close trajectories have the tendency to diverge is a typical property of a chaotic behaviour (it means that some Lyapounov exponent are positive).

## 5. CONCLUSION

For diluted networks (Sections 3 and 4), it is possible to extend some of the above calculations to describe more complex systems<sup>19-22</sup>. For example, one can produce short and long term memory effects by considering that the synapses  $J_{ij}$  are bounded ( $|J_{ij}| < L$ ) and that adding a new pattern changes the synapsis only if the constraint  $|J_{ij}| < L$  is satisfied<sup>19</sup>. One can also choose the  $J_{ij}$  in order to produce temporal sequences of patterns<sup>20</sup>.

The two models described in sections 2 and 3 have the following two simplifying features :

- (1) - there is no architecture : all the neurons play similar roles
- (2) - the synapses are given explicitly in terms of patterns

Recently it has been shown that none of these two simplifying assumptions is essential for the model to be soluble.

One can construct layered networks for which the dynamics can still be solved exactly<sup>4,23</sup>. The solution and the properties are similar to those

of the diluted model.

One can also use  $J_{ij}$  which are no longer given explicitly in terms of the patterns (like in the Hebb rule (7)) but which are arbitrary with the condition that there are attractors near the stored patterns<sup>24-27</sup>. This might be a first step in understanding the various learning rules which have been proposed and consist in using iterative procedures to modify the  $J_{ij}$  in order to create or to enlarge the basins of attraction near stored patterns<sup>28-30</sup>.

The first part of this short review has already appeared in<sup>31</sup>.

#### References

1. E. Bienenstock, F. Fogelman Soulie and G. Weisbuch eds : Disordered Systems and Biological Organisation, Springer Verlag, Heidelberg 1986
2. P. Peretto and J.J. Niez, in ref.[1] and Biol. Cybern. 54, 53 (1986)
3. T. Kohonen, Self Organization and Associative Memory Springer Verlag Berlin 1984
4. D.J. Amit, H. Gutfreund, H. Sompolinsky, Ann. Phys. 173, 30 (1987)
5. E. Domany, J. Stat. Phys. 51, 743 (1988)
6. P. Rujan, Worshop on "Systems with learning and memory abilities" 1987
7. W. Kinzel, preprint 1988 to appear in Physica Scripta
8. J.A. Hertz, preprint Nordita 88/25S
9. C. Campbell, D. Sherrington and K.Y.M. Wong, preprint 88.
10. J.P. Nadal, G. Toulouse, J.P. Changeux and S. Dehaene, Europhys. Lett. 1, 535 (1986)
11. M. Mezard, J.P. Nadal, G. Toulouse, J. Phys. Paris 47, 1457 (1986)
12. W.A. Little, Math. Biosci. 19, 101 (1974))

13. J.J. Hopfield, Proc. Nat. Acad. Sci. USA 79, 2554 (1982)
14. see for example K. Binder and A.P. Young, Rev. Mod. Phys. 58, 801 (1986) for a review on spin glasses)
15. B. Derrida : "Dynamics of Automata, Spin Glasses and Neural Network models" lectures given at Noto (Sicily) 1987 and references therein
16. D.J. Amit, H. Gutfreund, H. Sompolinsky, Phys. Rev. Lett. 55, 1530 (1985), Phys. Rev. A32, 1007 (1985)
17. B. Derrida, E. Gardner, A. Zippelius, Europhys. Lett. 4, 167 (1987)
18. B. Derrida, Preprint 88.
19. B. Derrida, J.P. Nadal, J. Stat. Phys. 49, 993 (1987)
20. H. Gutfreund, M. Mezard, preprint 87, see also J. Buhman, K. Schulten Europhys. Lett. 4, 1205 (1987)
21. A.J. Noest, Europhys. Lett. 6, 469 (1988)
22. R. Kree, A. Zippelius, Phys. Rev. A36, 4421 (1987).
23. R. Meir, E. Domany, Phys. Rev. Lett. 59, 359 (1987), Europhys. Lett. 4, 465 (1987), Phys. Rev. A37, 608 (1988)
24. E. Gardner, Europhys. Lett. 4, 1205 (1987)
25. E. Gardner, J. Phys. A21, 257 (1988)
26. E. Gardner, B. Derrida, J. Phys. A21, 271 (1988)
27. W. Krauth, M. Mezard, J. Phys. A 20, L745 (1987)
28. S. Diederich, M. Oppen, Phys. Rev. Lett. 58, 949 (1987)

29. E. Gardner, N. Stroud, D.J. Wallace, preprint 87
30. G. Pöppel, U. Krey, Europhys. Lett. 4, 979 (1987)
31. B. Derrida, Nucl. Phys. B (Proc. Supp.) 4, 673 (1988)