# Neural-network model composed of multidimensional spin neurons

Y. Nakamura, K. Torii, and T. Munakata

*Department of Applied Mathematics and Physics, Kyoto University, Kyoto 606-01 Japan*

(Received 15 August 1994)

As an extension of the Hopfield model, we propose a neural network composed of $D$-dimensional spin neurons ($D \geq 1$). Our model is equivalent to the Hopfield model in the case of $D = 1$ and is related to the clock neural network in the case of $D = 2$. We derive the free energy of our model using the replica symmetric theory. When a finite number of patterns are embedded they are found to be retrievable if the temperature $T$ is lower than $1/D$. The phase diagram and the storage capacity of the network are also obtained with the storage capacity $\alpha_c = 0.0743$ ($D = 2$) and $\alpha_c = 0.0432$ ($D = 3$) for $T = 0$.

PACS number(s): 87.10.+e, 05.20.−y

## I. INTRODUCTION

In the last decade, there has been a great deal of research on neural networks as associative memories. Various properties, for example the storage capacity of the network, were studied with use of techniques of statistical mechanics of spin systems [1–5]. Most of these studies, however, treated Ising type neurons which represent only two states, the firing and the resting one.

Recently, there is a growing interest in neural networks with multistate neurons. The merit of the model is that one neuron can express a complex state such as a color or shade of gray of each pixel in the pattern, which would otherwise require multiple Ising type neurons per pixel.

Rieger [6] and Bollé *et al.* [7,8] used neurons represented by spin variables $S_i$ ($i = 1, \ldots, N$), which can take $Q$ values,

$$S_i = -1 + \frac{2(k-1)}{Q-1} \qquad (i = 1, \ldots, N; k = 1, \ldots, Q).$$

$$\tag{1}$$

Earlier, Kanter proposed a model composed of Potts neurons with $q$ possible discrete states [9]. The dynamics of the Potts neural network is very complicated (see Refs. [9,10]).

Another possibility for a multistate neuron is the so-called circular representation, in which the state of the neuron is represented by points on the circle. The property of this model is that the state of each neuron can be expressed by a phase variable. We show here two examples of such models; one is the phasor model proposed by Noest [11,12] and the other is the clock model proposed by Cook [13]. In Sec. II, we will discuss how these models are similar to and different from our model.

Noest discussed a phasor network composed of unit-length two-dimensional vectors (phasors) as neurons. The network has $N$ phasors $S_i$ ($i = 1, \ldots, N$), which are complex numbers with $|S_i| = 1$. In the first version of the model [11], the $S_i$ could assume any such value, but in the second version [12], the $S_i$ were further restricted to the $q$ roots of $S_i^q \equiv 1$. When the $i$th phasor in the $\mu$th

pattern is described as $\xi_i^\mu$ ($i = 1, \ldots, N; \mu = 1, \ldots, p$), synaptic couplings $C_{ij}$ are defined as

$$C_{ij} = \frac{1}{N} \sum_{\mu=1}^{p} \xi_i^\mu \bar{\xi}_j^\mu (1 - \delta_{ij}),$$

$$\tag{2}$$

in order to store patterns. Note that $\delta_{ij}$ is Kronecker's delta function and $\bar{x}$ represents the complex conjugate of a complex $x$. Since $\xi_j^\mu$ can be expressed as $e^{i\theta_j^\mu}$, storing patterns is equivalent to storing phases. The dynamics of each phasors depends on the local (complex) field,

$$h_i = \sum_j C_{ij} S_j.$$

$$\tag{3}$$

For a discrete-time updating the state of $i$th phasor at the next step is defined as

$$S_i(t + \delta t) = \frac{h_i}{|h_i|},$$

$$\tag{4}$$

and for a continuous-time updating the $i$th phasor evolves as

$$\frac{dS_i}{dt} = h_i - \bar{h}_i S_i^2.$$

$$\tag{5}$$

In the network with an asynchronous discrete-time or a continuous-time updating, an energy can be defined as

$$E = -\frac{1}{2} \sum_{i \neq j} \bar{S}_i C_{ij} S_j,$$

$$\tag{6}$$

which plays the role of a Lyapunov function. Noest studied these networks with diluted synaptic connections. The case of full connection was studied by Gerl *et al.* [14].

Cook investigated a related $Q$-state clock neural-network model. The Hamiltonian of the system is

$$H = -\frac{1}{2N} \sum_{i \neq j} \sum_\mu \cos \frac{2\pi}{Q} \left\{ (n_i - \xi_i^\mu) - (n_j - \xi_j^\mu) \right\},$$

$$\tag{7}$$

where $N$ is the number of neurons, $n_i$ ($= 0, 1, \ldots, Q-$

1) $(i = 1, \ldots, N)$ the state of $i$th neuron and $\xi_i^\mu$ $(= 0, 1, \ldots, Q-1)$ $(i = 1, \ldots, N; \mu = 1, \ldots, p)$ the state of $i$th neuron in the $\mu$th pattern stored in the network. This system is reduced to the Hopfield model [15] for $Q = 2$ and to $x$-$y$ spin systems for $Q = \infty$, where the state becomes continuous. It is to be noted that the Hamiltonian (7), as well as the energy (6), is invariant under the transformation $\{n_i\} \rightarrow \{n_i + k \pmod{Q}\}$, where $k$ is an integer. Hence, if a pattern, say $\{\eta_i\}$, is stored, the $Q-1$ related configurations, $\{\eta_{k_i}\}$, are also stored, where

$$\eta_{k_i} = \eta_i + k \mod Q \quad (i = 1, \ldots, N; k = 1, \ldots, Q-1). \tag{8}$$

Therefore, notice that only the set of phase differences is meaningful as the information to be stored. Cook analyzed this model with the replica symmetric theory and estimated the storage capacity as $\alpha_c = 0.038$ in the limit $Q = \infty$.

Noest's model and Cook's model can be regarded as an extension of the Hopfield model to the clock type one. We can see similarities between the Noest's model and the Cook's model. They consist of essentially the same type of model neurons: continuous versions of the Noest's model correspond to the Cook's model with infinite $Q$ and discrete versions of the Noest's model correspond to the Cook's model with finite $Q$. Furthermore, they have the same energy as we shall see in the next section. In this paper we want to consider a different, and more general extension of the Hopfield model using $D$-dimensional spins (i.e., $D$-dimensional unit vectors) as neurons. Our model is a natural $D$-dimensional $(D = 1, 2, \ldots)$ generalization of the Hopfield model. For $D = 2$, our model is similar to those of Noest and Cook, but the energy of our model, or more precisely the algebraic structure of the interactions, is quite different from the models of Noest and Cook. These distinctions will be discussed in detail in the next section.

This paper is organized as follows: In Sec. II, we propose our neural-network model composed of multi-dimensional spin neurons. Simulation results are shown in Sec. III. In Sec. IV, we analyze our model theoretically using a replica symmetric theory. We derive the free energy of our model near saturation. On the basis of these results, we calculate and discuss the phase diagram and the storage capacity. In the last section, we summarize our results.

## II. MODEL

We consider a network composed of $N$ neurons, which are described by $D$-dimensional unit vectors $\boldsymbol{x}_i = {}^t(x_{i(1)}, x_{i(2)}, \ldots, x_{i(D)})$ $(i = 1, \ldots, N)$, where ${}^t\boldsymbol{z}$ denotes the transposed vector of $\boldsymbol{z}$. Each neuron represents an arbitrary point on a surface of $D$-dimensional unit sphere. Especially in the case of $D = 2$, this model neuron is similar to that of Cook $(Q = \infty)$ and of Noest. However, the neurons are interconnected with all the others through synaptic couplings, which differ fundamentally from those in the Noest and Cook models. In this paper,

we define synaptic couplings by an extended Hebb's rule as

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^{p} \boldsymbol{\xi}_i^{\mu t} \boldsymbol{\xi}_j^\mu \quad (i \neq j), \qquad J_{ii} = O, \tag{9}$$

where $\boldsymbol{\xi}_i^\mu$ denotes the state of $i$th neuron in the $\mu$th memorized pattern $(\mu = 1, \ldots, p)$ and is also a $D$-dimensional unit vector. Note that a synaptic strength $J_{ij}$ is here a $D \times D$ matrix and $O$ is a $D \times D$ zero matrix.

There are two kinds of dynamics, one is under the zero temperature $(T = 0)$ and the other under the finite temperature $(T \neq 0)$. First, we consider the case $T = 0$. We calculate the local field $\boldsymbol{h}_i(t)$ on the $i$th neuron at time $t$ as

$$\boldsymbol{h}_i(t) = \sum_{j=1}^{N} J_{ij} \boldsymbol{x}_j(t). \tag{10}$$

$\boldsymbol{h}_i(t)$ is a $D$-dimensional and, in general, a nonunit vector. The state of the neuron $i$ at the next time step is determined by normalizing $\boldsymbol{h}_i(t)$ as

$$\boldsymbol{x}_i(t + \Delta t) = \frac{\boldsymbol{h}_i(t)}{|\boldsymbol{h}_i(t)|}. \tag{11}$$

As to the way of updating, here we adopt an asynchronous one, for which we can define an energy function,

$$E = -\frac{1}{2} \sum_{i \neq j} {}^t\boldsymbol{x}_i J_{ij} \boldsymbol{x}_j. \tag{12}$$

This energy is a Lyapunov function for the dynamics. This is seen as follows: Suppose one of the neuron, say $i$th neuron, updates as $\boldsymbol{x}_i(t) \rightarrow \boldsymbol{x}_i(t+\Delta t)$, then the variation of energy $\Delta E$ is calculated as

$$\Delta E = -|\boldsymbol{h}_i(t)|\{1 - {}^t\boldsymbol{x}_i(t)\, \boldsymbol{x}_i(t + \Delta t)\} \leq 0. \tag{13}$$

Next, in the case of $T \neq 0$, a random contribution is added to Eq. (10) so that the equilibrium state of the system realizes the canonical ensemble [5,11].

We now discuss the differences between our model and those of Noest and Cook. For comparison we restrict our model to the two-dimensional case (clock type). In the Noest's model, when $S_j$ is described as $e^{i\varphi_j}$ and $\xi_j^\mu$ as $e^{i\theta_j^\mu}$, the energy (6) is rewritten as

$$E = -\frac{1}{2N} \sum_{i \neq j} \sum_{\mu} \cos\{(\varphi_i - \theta_i^\mu) - (\varphi_j - \theta_j^\mu)\} . \tag{14}$$

In the Cook's model, with new notations $(2\pi/Q)n_i = \varphi_i$ and $(2\pi/Q)\xi_i^\mu = \theta_i^\mu$, the expression (7) coincides with Eq. (14). When we describe $\boldsymbol{x}_i = {}^t(\cos \varphi_i, \ \sin \varphi_i)$ and $\boldsymbol{\xi}_i^\mu = {}^t(\cos \theta_i^\mu, \ \sin \theta_i^\mu)$, the energy of our model (12) is rewritten as

$$E = -\frac{1}{2N} \sum_{i \neq j} \sum_{\mu} \cos(\varphi_i - \theta_i^\mu) \cos(\varphi_j - \theta_j^\mu). \tag{15}$$

Therefore, the translational symmetry in Noest's and

Cook's model is not existent in our model, thus the translated phase pattern $\{\theta_i + \bar{\theta}\}$ is quite different from the pattern $\{\theta_i\}$ in our model. However, our model has an inversion symmetry in which the energy (12) is invariant under the transformation $\{\boldsymbol{x}_i\} \to \{-\boldsymbol{x}_i\}$ or $\{\boldsymbol{\xi}_i^\mu\} \to \{-\boldsymbol{\xi}_i^\mu\}$. This fact means that inverted patterns can also be memorized as in the Hopfield model.

## III. SIMULATION

In this section, we simulate our model according to the dynamics (10) and (11) with synaptic couplings (9) under the zero temperature. We adopt an asynchronous updating in which at each time step, one neuron is selected randomly to be updated. The time step is set as $\Delta t = 1/N$. Simulation results are shown in the following.

First we examine the case $D = 2$. Patterns are written as $\boldsymbol{\xi}_i^\mu = {}^t(\cos\theta_i^\mu, \ \sin\theta_i^\mu)$, with $\theta_i^\mu$ distributed uniformly in $[0, 2\pi]$. These patterns are uncorrelated in the sense of $\sum_{i=1}^N \langle \boldsymbol{\xi}_i^\mu \boldsymbol{\xi}_i^\nu \rangle = \delta_{\mu\nu}$. Figure 1 shows the time evolutions of overlap $m^1(t)$ defined as

$$m^1(t) = \frac{1}{N}\sum_{j=1}^N {}^t\boldsymbol{\xi}_i^1 \, \boldsymbol{x}_i, \qquad (16)$$

which start from several initial value of $m^1(0)$. From Fig. 1(a), which is for the case of $\alpha \equiv p/N = 0.05$ ($p = 20, N = 400$), we see that the network can retrieve a memorized pattern. On the other hand, in the case of $\alpha = 0.1$ ($p = 40, N = 400$) the network cannot retrieve a memorized pattern [Fig. 1(b)]. These results tell us that the storage capacity $\alpha_c$ is $0.05 < \alpha_c < 0.1$. Note that this storage capacity is larger than that of Cook ($\alpha_c = 0.038$ for $Q = \infty$) [13]. This will be discussed theoretically in the following section.

Next, we make simulations in the case of $D = 3$. In this case, $p$ patterns are expressed as $\boldsymbol{\xi}_i^\mu = {}^t(\sin\psi_i^\mu \cos\theta_i^\mu, \ \sin\psi_i^\mu \sin\theta_i^\mu, \ \cos\psi_i^\mu)$, with uniformly distributed $\psi_i^\mu$ in $[0, \pi]$ and $\theta_i^\mu$ in $[0, 2\pi]$. These patterns are also uncorrelated. The result is that, in the case of $\alpha = 0.025 (p = 10, N = 400)$, the network can retrieve a memorized pattern, but in the case of $\alpha = 0.075$ ($p = 30, N = 400$), the network cannot retrieve a memorized pattern. Therefore, the storage capacity $\alpha_c$ is $0.025 < \alpha_c < 0.075$. This result will also be considered theoretically in the following section.

## IV. MEAN-FIELD THEORY

The mean-field theory is performed to calculate the free energy of the system with use of the replica method, i.e., the free energy per neuron, $f$, is written as

$$f = \lim_{n \to 0} \lim_{N \to \infty} -\frac{\langle\langle Z^n \rangle\rangle - 1}{nN\beta}, \qquad (17)$$
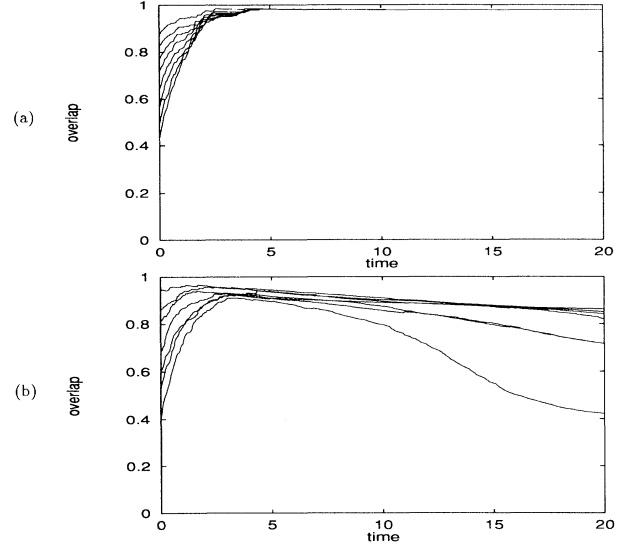


FIG. 1. Time evolutions of overlaps which start from several initial conditions. A time unit is $N$ steps. (a) $N = 400$, $p = 20$, $\alpha = 0.05$; (b) $N = 400$, $p = 40$, $\alpha = 0.1$.

where $\langle\langle \ \rangle\rangle$ denotes a quenched average over the memorized patterns, $\{\boldsymbol{\xi}_i^\mu\}$, $\beta = 1/T$ an inverse temperature of the system, $Z$ the partition function of the system, and $n$ the number of replicas. In the following calculation we use the framework of Amit, Gutfreund, and Sompolinsky [3] (hereafter referred to as AGS). We adopt

$$H^\rho = -\frac{1}{2}\sum_{i \neq j} {}^t\boldsymbol{x}_i^\rho J_{ij}\boldsymbol{x}_j^\rho - \sum_{\nu=1}^s h^\nu \sum_{i=1}^N {}^t\boldsymbol{\xi}_i^\nu \, \boldsymbol{x}_i^\rho \qquad (18)$$

as the Hamiltonian of the $\rho$th replica of the system. We assume that there are a finite number $s$ ($\ll p$) of fields $h^\nu$ ($\nu = 1, \ldots, s$) each of which is coupled to $\nu$th condensed pattern. Then $\langle\langle Z^n \rangle\rangle$ is calculated as

$$\langle\langle Z^n \rangle\rangle = \left\langle\!\!\left\langle \mathrm{Tr}_{\boldsymbol{x}^\rho} \exp\left( -\beta \sum_{\rho=1}^n H^\rho \right) \right\rangle\!\!\right\rangle \qquad (19)$$

$$= \left\langle\!\!\left\langle \mathrm{Tr}_{\boldsymbol{x}^\rho} \exp\left[ \frac{\beta}{2N}\sum_{ij\mu\rho} ({}^t\boldsymbol{\xi}_i^\mu \, \boldsymbol{x}_i^\rho)({}^t\boldsymbol{\xi}_j^\mu \, \boldsymbol{x}_j^\rho) \right.\right.$$

$$\left.\left. -\frac{\beta}{2N}\sum_{i\mu\rho} ({}^t\boldsymbol{\xi}_i^\mu \, \boldsymbol{x}_i^\rho)^2 + \beta\sum_{\nu\rho} h^\nu \sum_i {}^t\boldsymbol{\xi}_i^\nu \, \boldsymbol{x}_i^\rho \right] \right\rangle\!\!\right\rangle.$$

$$(20)$$

Using the Stratonovich-Hubbard transformation,

$$\exp\left( \frac{A}{2}a^2 \right) = \frac{1}{\sqrt{2\pi A}}\int_{-\infty}^\infty dx \exp\left( -\frac{x^2}{2A} + ax \right), \quad (21)$$

we can proceed to calculate Eq. (20) as

$$\langle\!\langle Z^n \rangle\!\rangle = (\beta N)^{\frac{pn}{2}} \left\langle\!\!\left\langle \text{Tr}_{x^\rho} \int \prod_{\mu\rho} \frac{dm_\rho^\mu}{\sqrt{2\pi}} \exp\left[\beta N \left\{ -\frac{1}{2}\sum_{\mu\rho}(m_\rho^\mu)^2 + \sum_{\mu\rho} m_\rho^\mu \left(\frac{1}{N}\sum_i {}^t\xi_i^\mu x_i^\rho\right) - \sum_{\mu\rho}\frac{1}{2N^2}\sum_i ({}^t\xi_i^\mu x_i^\rho)^2 \right\}\right]\right.$$

$$\left.\times \int \prod_{\nu\rho} \frac{dm_\rho^\nu}{\sqrt{2\pi}} \exp\left[\beta N \left\{ -\frac{1}{2}\sum_{\nu\rho}(m_\rho^\nu)^2 + \sum_{\nu\rho}(m_\rho^\nu + h^\nu) \left(\frac{1}{N}\sum_i {}^t\xi_i^\nu x_i^\rho\right) - \sum_{\nu\rho}\frac{1}{2N^2}\sum_i ({}^t\xi_i^\nu x_i^\rho)^2 \right\}\right]\right\rangle\!\!\right\rangle. \quad (22)$$

The sums $\sum_\nu$ and $\sum_\mu$ are over the first $s$ patterns and over the remaining $p-s$ patterns, respectively. The products $\prod_\nu$ and $\prod_\mu$ are also considered in the same way. In order to average the first exponential in Eq. (22) over the $p-s$ uncondensed patterns $\{\xi_i^\mu\}$, we make variable transformations,

$$x_i^\rho = \begin{pmatrix} x_{i(1)}^\rho \\ x_{i(2)}^\rho \\ x_{i(3)}^\rho \\ \vdots \\ x_{i(D)}^\rho \end{pmatrix} = \begin{pmatrix} \sin\zeta_{i(1)}^\rho \sin\zeta_{i(2)}^\rho \cdots \sin\zeta_{i(D-2)}^\rho \cos\varphi_i^\rho \\ \sin\zeta_{i(1)}^\rho \sin\zeta_{i(2)}^\rho \cdots \sin\zeta_{i(D-2)}^\rho \sin\varphi_i^\rho \\ \sin\zeta_{i(1)}^\rho \cdots\cdots\cdots \cos\zeta_{i(D-2)}^\rho \\ \vdots \\ \cos\zeta_{i(1)}^\rho \end{pmatrix}, \quad (23)$$

and

$$\xi_i^\mu = \begin{pmatrix} \xi_{i(1)}^\mu \\ \xi_{i(2)}^\mu \\ \xi_{i(3)}^\mu \\ \vdots \\ \xi_{i(D)}^\mu \end{pmatrix} = \begin{pmatrix} \sin\psi_{i(1)}^\mu \sin\psi_{i(2)}^\mu \cdots \sin\psi_{i(D-2)}^\mu \cos\theta_i^\mu \\ \sin\psi_{i(1)}^\mu \sin\psi_{i(2)}^\mu \cdots \sin\psi_{i(D-2)}^\mu \sin\theta_i^\mu \\ \sin\psi_{i(1)}^\mu \cdots\cdots\cdots \cos\psi_{i(D-2)}^\mu \\ \vdots \\ \cos\psi_{i(1)}^\mu \end{pmatrix}, \quad (24)$$

using polar coordinates in a $D$-dimensional space. $\zeta_{i(k)}^\rho, \psi_{i(k)}^\mu$ $(k=1,\ldots,D-2)$ vary from 0 to $\pi$ and $\varphi_i^\rho, \theta_i^\mu$ from 0 to $2\pi$. With use of these transformation, we can calculate an average over $\{\xi_i^\mu\}$ as

$$\langle\!\langle f(\xi_i^\mu) \rangle\!\rangle_{\xi_i^\mu} = \frac{1}{\mathcal{J}} \int_0^\pi d\psi_{i(1)}^\mu \int_0^\pi d\psi_{i(2)}^\mu \cdots \int_0^\pi d\psi_{i(D-2)}^\mu \int_0^{2\pi} d\theta_i^\mu \sin^{(D-2)}\psi_{i(1)}^\mu \sin^{(D-3)}\psi_{i(2)}^\mu \cdots \sin\psi_{i(D-2)}^\mu$$
$$\times f(\psi_{i(1)}^\mu, \ldots, \psi_{i(D-2)}^\mu, \theta_i^\mu) \quad (25)$$

where

$$\mathcal{J} = \int_0^\pi d\psi_{i(1)}^\mu \int_0^\pi d\psi_{i(2)}^\mu \cdots \int_0^\pi d\psi_{i(D-2)}^\mu \int_0^{2\pi} d\theta_i^\mu. \quad (26)$$

Therefore, the average is calculated as

$$\left\langle\!\!\left\langle \exp\left[\beta N \left\{ \sum_{\mu\rho} m_\rho^\mu \left(\frac{1}{N}\sum_i {}^t\xi_i^\mu x_i^\rho\right) - \sum_{\mu\rho}\frac{1}{2N^2}\sum_i ({}^t\xi_i^\mu x_i^\rho)^2 \right\}\right]\right\rangle\!\!\right\rangle_{\xi_i^\mu}$$

$$= \prod_{\mu i} \left\langle\!\!\left\langle \exp\left[\beta \sum_\rho \left\{ m_\rho^{\mu\,t}\xi_i^\mu x_i^\rho - \frac{1}{2N}({}^t\xi_i^\mu x_i^\rho)^2 \right\}\right]\right\rangle\!\!\right\rangle_{\xi_i^\mu} = \exp\left(-\frac{\beta pn}{2D}\right) \exp\left[\sum_{\mu i}\frac{\beta^2}{2D}\sum_{\rho\sigma} m_\rho^\mu m_\sigma^{\mu\,t} x_i^\rho x_i^\sigma\right] \quad (27)$$

(see Appendix). Inserting Eq. (27) into Eq. (22) and introducing order parameters $r_{\rho\sigma}$ and $q_{\rho\sigma}$, we integrate over $m_\rho^\mu$ and make a saddle-point approximation for the integration over $m_\rho^\nu, r_{\rho\sigma}$, and $q_{\rho\sigma}$, following AGS. Thus, we get the expression of $\langle\!\langle Z^n \rangle\!\rangle$ as

$$\langle\!\langle Z^n \rangle\!\rangle = \exp\left(-\frac{\beta pn}{2D}\right) \left(\frac{\alpha\beta^2}{2D}\right)^{\frac{n(n-1)}{2}} \exp\left[-\frac{p}{2}\text{Tr}\ln\left\{\left(1-\frac{\beta}{D}\right)I - \frac{\beta}{D}Q\right\} - \frac{N\alpha\beta^2}{2D}\sum_{\rho\neq\sigma} r_{\rho\sigma}q_{\rho\sigma}\right.$$

$$\left. -\frac{\beta N}{2}\sum_{\nu\rho}(m_\rho^\nu)^2 + \left\langle\!\!\left\langle N\ln\text{Tr}_{x^\rho}\exp\left\{\frac{\alpha\beta^2}{2D}\sum_{\rho\neq\sigma} r_{\rho\sigma}{}^t x^\rho x^\sigma + \beta\sum_{\nu\rho}(m_\rho^\nu + h^\nu){}^t\xi^\nu x^\rho - \frac{\beta}{2N}({}^t\xi^\nu x^\rho)^2\right\}\right\rangle\!\!\right\rangle_{\xi^\nu}\right], \quad (28)$$

where $I$ is a unit matrix with $n \times n$ elements and $Q$ is a matrix $\{q_{\rho\sigma}\}$ with zero diagonal elements. The parameters $m_\rho^\nu, q_{\rho\sigma}$, and $r_{\rho\sigma}$ are determined from the saddle-point equations to be given as follows:

$$m_\rho^\nu = \left\langle\!\!\left\langle \frac{1}{N} \sum_i {}^t\boldsymbol{\xi}_i^\nu \langle \boldsymbol{x}_i^\rho \rangle \right\rangle\!\!\right\rangle, \qquad \nu = 1, \dots, s, \qquad (29)$$

$$q_{\rho\sigma} = \left\langle\!\!\left\langle \frac{1}{N} \sum_i {}^t\langle \boldsymbol{x}_i^\rho \rangle \langle \boldsymbol{x}_i^\sigma \rangle \right\rangle\!\!\right\rangle, \qquad (30)$$

$$r_{\rho\sigma} = \frac{1}{\alpha} \sum_{\mu > s} \left\langle\!\!\left\langle m_\rho^\mu m_\sigma^\mu \right\rangle\!\!\right\rangle, \qquad (31)$$

where $\langle\,\rangle$ denotes a thermal average. Eventually, we get the free energy per neuron as

$$f = \lim_{n \to 0} \left[ \frac{\alpha}{2D} + \frac{\alpha}{2\beta n} \operatorname{Tr} \ln \left\{ \left( 1 - \frac{\beta}{D} \right) I - \frac{\beta}{D} Q \right\} \right.$$

$$+ \frac{\alpha\beta}{2Dn} \sum_{\rho \neq \sigma} r_{\rho\sigma} q_{\rho\sigma} + \frac{1}{2n} \sum_{\nu\rho} (m_\rho^\nu)^2$$

$$\left. - \frac{1}{n\beta} \left\langle \ln \operatorname{Tr}_{\boldsymbol{x}^\rho} \exp(\beta H_\xi) \right\rangle_{\xi^\nu} \right], \qquad (32)$$

$$H_\xi = \frac{\alpha\beta}{2D} \sum_{\rho \neq \sigma} r_{\rho\sigma} {}^t\boldsymbol{x}^\rho \boldsymbol{x}^\sigma + \sum_{\nu\rho} (m_\rho^\nu + h^\nu) {}^t\boldsymbol{\xi}^\nu \boldsymbol{x}^\rho. \qquad (33)$$

In the case of $D = 1$, this result exactly reduces to that of AGS for the Hopfield model [3,4]. When we take the replica symmetric assumption in which $m_\rho^\nu = m^\nu, q_{\rho\sigma} = q$, and $r_{\rho\sigma} = r$, Eq. (32) becomes

$$f = \frac{\alpha}{2D} + \frac{\alpha}{2\beta} \ln \left\{ 1 - \frac{\beta}{D}(1 - q) \right\} - \frac{\alpha}{2} \frac{q}{D - \beta(1 - q)}$$

$$+ \frac{\alpha\beta r(1 - q)}{2D} + \frac{1}{2} \sum_\nu (m^\nu)^2$$

$$- \frac{1}{\beta} \left\langle\!\!\left\langle \ln \operatorname{Tr}_{\boldsymbol{x}} \exp \left( \sum_k A_k x_{(k)} \right) \right\rangle\!\!\right\rangle, \qquad (34)$$

$$A_k = \beta \left\{ \sqrt{\frac{\alpha r}{D}} z_k + \sum_\nu (m^\nu + h^\nu) \xi_{(k)}^\nu \right\} \qquad (k = 1, \dots, D), \qquad (35)$$

where $\langle\!\langle\,\rangle\!\rangle$ denotes the combined average over condensed patterns $\{\boldsymbol{\xi}^\nu\}$ [Eq. (25)] and over the Gaussian noises $z_1, \dots, z_k$ defined as

$$\langle\!\langle f(\boldsymbol{\xi}^\nu) \rangle\!\rangle = \left\langle\!\!\left\langle \int_{-\infty}^\infty \cdots \int_{-\infty}^\infty \frac{dz_1 \cdots dz_D}{(2\pi)^{D/2}} \right.\right.$$

$$\left.\left. \times \exp \left( -\frac{z_1^2 + \cdots + z_D^2}{2} \right) f(\boldsymbol{\xi}^\nu) \right\rangle\!\!\right\rangle_{\xi^\nu}. \qquad (36)$$

Mean-field equations are determined as follows:

$$\frac{\partial f}{\partial m^\nu} = 0,$$

$$m^\nu = \frac{1}{\beta} \left\langle\!\!\left\langle \frac{\partial}{\partial m^\nu} \left[ \ln \operatorname{Tr}_{\boldsymbol{x}} \exp \left( \sum_k A_k x_{(k)} \right) \right] \right\rangle\!\!\right\rangle$$

$$\frac{\partial f}{\partial r} = 0, \qquad (37)$$

$$\beta(1 - q) = \frac{2D}{\alpha\beta} \left\langle\!\!\left\langle \frac{\partial}{\partial r} \left[ \ln \operatorname{Tr}_{\boldsymbol{x}} \exp \left( \sum_k A_k x_{(k)} \right) \right] \right\rangle\!\!\right\rangle$$

$$\frac{\partial f}{\partial q} = 0, \qquad (38)$$

$$r = \frac{Dq}{\{D - \beta(1 - q)\}^2}. \qquad (39)$$

We note that when we solve these equations, a set of solutions $m^\mu = 0, q = 0$ represents the paramagnetic state, $m^\mu = 0, q \neq 0$ the spin-glass state and $m^\mu \neq 0, q \neq 0$ the retrieval state. For the case of $D = 2$, $\operatorname{Tr}_{\boldsymbol{x}}$ in Eq. (34) is calculated explicitly to give

$$\operatorname{Tr}_{\boldsymbol{x}} \exp \left( \sum_k A_k \boldsymbol{x}_{(k)} \right)$$

$$= \int_0^{2\pi} d\varphi \exp(A_1 \cos\varphi + A_2 \sin\varphi)$$

$$= 2\pi I_0 \left( \sqrt{A_1^2 + A_2^2} \right), \qquad (40)$$

where $I_k(z)$ is the $k$th order modified Bessel function defined by

$$I_k(z) = \frac{1}{2\pi} \int_0^{2\pi} d\phi\, e^{z\cos\phi} \cos k\phi. \qquad (41)$$

In the following subsections, we examine the properties of our model using the free energy (34) under the conditions

$$m^1 = m, \quad m^\nu = 0 \ (\nu \geq 2), \quad h^\nu = 0. \qquad (42)$$

### A. $\alpha = 0$

In this subsection we examine the case $\alpha = 0$, in which an intensive number of patterns are embedded in the network ($N = \infty$). We put $\alpha = 0$ in Eq. (34) and get

$$f = \frac{1}{2} m^2 - \frac{1}{\beta} \ln C_D \frac{I_{\frac{D}{2}-1}(\beta m)}{(\frac{1}{2}\beta m)^{\frac{D}{2}-1}}, \qquad (43)$$

where

$$C_1 = 2\sqrt{\pi}, \quad C_2 = 2\pi, \quad C_3 = 2\pi\sqrt{\pi},$$

$$C_D = \sqrt{\pi}\Gamma \left( \frac{D-1}{2} \right) \int_0^\pi d\zeta_2 \sin^{D-3}\zeta_2$$

$$\times \cdots \int_0^\pi d\zeta_{D-2} \sin\zeta_{D-2} \int_0^{2\pi} d\varphi \quad (D \geq 4).$$

The mean-field equation becomes

$$m = \frac{I_{\frac{D}{2}}(\beta m)}{I_{\frac{D}{2}-1}(\beta m)} \equiv g(\beta; m). \qquad (44)$$

According to the property of the modified Bessel function,

$$\frac{d}{dz}\frac{I_\nu(z)}{I_{\nu-1}(z)} > 0, \quad \frac{d^2}{dz^2}\frac{I_\nu(z)}{I_{\nu-1}(z)} \begin{cases} < 0 & (z > 0) \\ > 0 & (z < 0) \end{cases}, \qquad (45)$$

the function $g(\beta; m)$ is a sigmoid type one. Therefore, under the condition

$$\left.\frac{d}{dm}g(\beta; m)\right|_{m=0} \leq 1, \qquad (46)$$

which yields $\beta \leq D$ or $T \geq 1/D$, Eq. (44) has no solution except for $m = 0$. On the other hand, as the noise $T$ decreases from $1/D$ or $\beta$ increases from $D$, Eq. (44) begins to have a nonzero positive solution which is stable. Figure 2 shows positive solutions of Eq. (44) as a function of temperature $T = 1/\beta$, for various case of dimension $D$.

### B. $T = 0$

Next we consider the case of a low temperature limit $T = 0$ or $\beta = \infty$, where the mean-field equations (37), (38), and (39) become as follows:

$$m = \left\langle\!\!\left\langle \frac{\sum_k z_k\xi_{(k)} + y}{\sqrt{\sum_k z_k^2 + 2y\sum_k z_k\xi_{(k)} + y^2}} \right\rangle\!\!\right\rangle \equiv f_1(y), \qquad (47)$$

$$\beta(1-q) = \left\langle\!\!\left\langle \frac{\sum_k z_k^2 + y\sum_k z_k\xi_{(k)}}{\sqrt{\sum_k z_k^2 + 2y\sum_k z_k\xi_{(k)} + y^2}} \right\rangle\!\!\right\rangle \equiv f_2(y), \qquad (48)$$

$$r = \frac{1}{D\left\{1 - \frac{1}{D}\sqrt{\frac{D}{\alpha r}}f_2(y)\right\}^2}, \qquad (49)$$
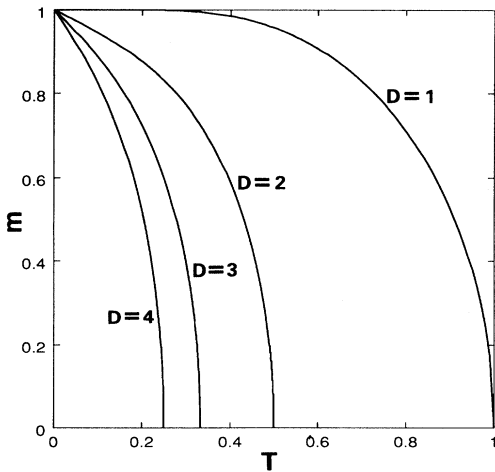
$$y \equiv \sqrt{\frac{D}{\alpha r}}m. \qquad (50)$$

These equations are reduced to a single equation for the variable $y$,

$$y = \frac{Df_1(y)}{\sqrt{\alpha} + f_2(y)}. \qquad (51)$$

This is a relation between the retrieval quality $m$ and the storage level $\alpha$. The storage capacity $\alpha_c$ is the value of $\alpha$, above which the Eq. (51) has no solution except for $y = 0$. The graphical solution of Eq. (51) is shown in Fig. 3 for $m \geq 0$. The straight line represents the left hand side (lhs). The dashed curves represent the right hand side (rhs) plotted for two values of $\alpha$, one below and one above $\alpha_c$ in the case of $D = 2$. For $\alpha < \alpha_c$, we have three non-negative solutions, $m_1 = 0, 0 < m_2 < m_3$. $m_1, m_3$ are stable and $m_2$ is unstable. The solution $m_1$ represents the spin-glass state because $m_1 = 0$ and $q \neq 0$ and $m_3$ the retrieval state because $m_1 \neq 0$ and $q \neq 0$. For $\alpha > \alpha_c$, there exists only one solution $m = 0$ which is the spin-glass state because of $q \neq 0$. Figure 3 tells us that the retrieval solution disappear abruptly. Figure 4 shows the solution $m$ as a function of $\alpha$ in three cases $D = 1, 2$, and 3. In the case of $D = 2$, retrieval solution disappears abruptly at $\alpha = 0.0743$, which is the storage capacity $\alpha_c$. Note that this storage capacity is larger than that of the Cook's model [13] and is consistent with the simulation results in Sec. II. In the case of $D = 3$, the storage capacity is calculated to be $\alpha_c = 0.0432$.

### C. $T$-$\alpha$ phase diagram

We now turn to the full mean-field equations (37), (38), and (39), keeping $T$ and $\alpha$ finite. In the case of $D = 2$, mean-field equations can be written explicitly as

$$m = \left\langle\!\!\left\langle \frac{\beta}{\sqrt{A_1^2 + A_2^2}} \frac{I_1\left(\sqrt{A_1^2 + A_2^2}\right)}{I_0\left(\sqrt{A_1^2 + A_2^2}\right)} \right.\right.$$
$$\left.\left. \times \left\{\sqrt{\frac{\alpha r}{2}}(z_1\xi_{(1)} + z_2\xi_{(2)}) + m\right\} \right\rangle\!\!\right\rangle, \qquad (52)$$



FIG. 2. Nonzero solutions of Eq. (44) as a function of temperature $T = 1/\beta$ for $D = 1, 2, 3$, and 4.
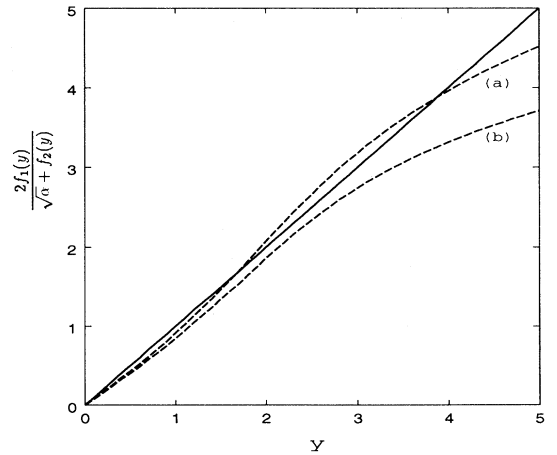


FIG. 3. The graphical representation of the solutions of Eq. (51) in the case $D = 2$. (a) $\alpha = 0.05$, (b) $\alpha = 0.1$.

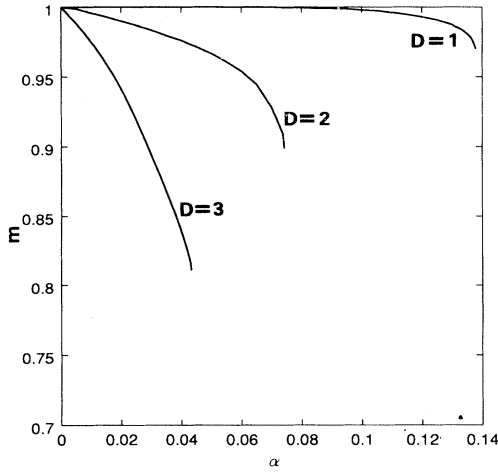FIG. 4. Solutions $m$ of Eqs. (50), (51) as a function of $\alpha$ in the case of $D = 1, 2$, and 3.

$$\beta(1 - q) = \sqrt{\frac{2}{\alpha r}} \left\langle\!\!\left\langle \frac{\beta}{\sqrt{A_1^2 + A_2^2}} \frac{I_1\left(\sqrt{A_1^2 + A_2^2}\right)}{I_0\left(\sqrt{A_1^2 + A_2^2}\right)} \right. \right.$$
$$\left. \left. \times \left\{\sqrt{\frac{\alpha r}{2}}(z_1^2 + z_2^2) + m(z_1\xi_{(1)} + z_2\xi_{(2)})\right\} \right\rangle\!\!\right\rangle,$$

$$\tag{53}$$

$$r = \frac{2q}{\{2 - \beta(1 - q)\}^2}. \tag{54}$$

The $T$-$\alpha$ phase diagram for $D = 2$ can be obtained by solving Eqs. (52), (53), and (54) numerically. However, because of the difficulty of the numerical calculation, we obtained only a part of the phase diagram, i.e., the boundary between the paramagnetic phase and the spin-glass one. In order to gain as much information about the phase diagram as possible under these circumstances, we try to obtain the qualitative structure of the phase diagram analytically.

At high temperature $T$ or low $\beta$, only the set of solutions $m = 0, q = 0$ (the paramagnetic state) is possible. Decreasing the temperature for fixed $\alpha$, we cross the transition temperature $T_g(\alpha)$ below which a set of solutions $m = 0, q \neq 0$ (the spin-glass state) appears. With the anticipation that $q$ will develop continuously from zero, the r.h.s. of Eq. (54) is expanded in powers of $q$ to give in the lowest order,

$$r \approx \frac{2q}{(2 - \beta)^2}. \tag{55}$$

Setting $m = 0$ in Eq. (53) and using Eq. (55), we get

$$q = \frac{\beta^2 \alpha r}{4} = \frac{\beta^2 \alpha q}{2(2 - \beta)^2}. \tag{56}$$

From this equation, we find the transition temperature $T_g(\alpha)$ to be

$$T_g(\alpha) = \frac{1}{2} + \frac{\sqrt{\alpha}}{2\sqrt{2}}. \tag{57}$$

When we, furthermore, decrease the temperature from $T_g(\alpha)$ for fixed $\alpha < \alpha_c = 0.0743$, we cross the transition temperature $T_M(\alpha)$. Below $T_M(\alpha)$, a set of solutions $m \neq 0, q \neq 0$ (the retrieval state) is possible. In this case the transition is a first order type, so there can be no expansion in $m$ generally. But we can make an analytic calculation in the corner of the phase diagram near $\alpha = 0$ and $T = \frac{1}{2}$, because we expect both $q$ and the discontinuity in $m$ to be small there. There are three small parameters $m, q$, and $t = \frac{1}{2} - T$. The equations (52), (53), and (54) are expanded in powers of these parameters to give

$$t = \frac{1}{2}m^2 + \alpha r, \tag{58}$$

$$q = 2m^2 + 2\alpha r, \tag{59}$$

$$r = \frac{q}{2(q - 2t)^2}. \tag{60}$$

From these equations we get

$$g(y) \equiv y^3 - 2\tau y^2 + y + 2\tau = 0 , \tag{61}$$

where

$$\tau \equiv \frac{t}{\sqrt{\alpha}}, \quad y \equiv \frac{m^2}{\sqrt{\alpha}}. \tag{62}$$

This equation has either two positive solutions or none at all. $T_M(\alpha)$ is determined by the disappearance of the two solutions. The value of $\tau$ at which the two solutions just disappear is calculated to be $\tau = 1.67$. Hence $T_M(\alpha)$ near $\alpha = 0$ and $T = 1/2$ is found to be

$$T_M(\alpha) = \frac{1}{2} - 1.67\sqrt{\alpha}. \tag{63}$$

In the case of $D = 3$, we can calculate $T_g(\alpha)$ and $T_M(\alpha)$ in the same way as shown above, to give

$$T_g(\alpha) = \frac{1}{3} + \frac{1}{3}\sqrt{\frac{5}{6}\alpha}, \tag{64}$$

$$T_M(\alpha) = \frac{1}{3} - 1.42\sqrt{\alpha}. \tag{65}$$

Summarizing the above results and the storage capacity, which is discussed in the previous subsection, we get a phase diagram as is shown in Fig. 5. We note that the agreement between the line $T_g$ obtained theoretically and that obtained numerically is excellent. Dashed lines are expected ones, which are depicted to guide the eyes. This figure tells us that the retrieval regime in the phase diagram becomes smaller as the dimension $D$ increases, as expected.

## V. SUMMARY

In this paper, we studied a neural network composed of $D$-dimensional spin neurons as an extension of the Hopfield model to a multidimensional one. In the case of $D = 2$, we found that a stored phase pattern $\{\theta_i\}$ itself, not a translated one $\{\theta_i + \bar{\theta}\}$, is retrieved. We analyzed the network by means of the replica symmetric theory and got the free energy of the network in the case of general dimension $D$. For $D = 1$, the free energy was confirmed to coincide that of the Hopfield model. With
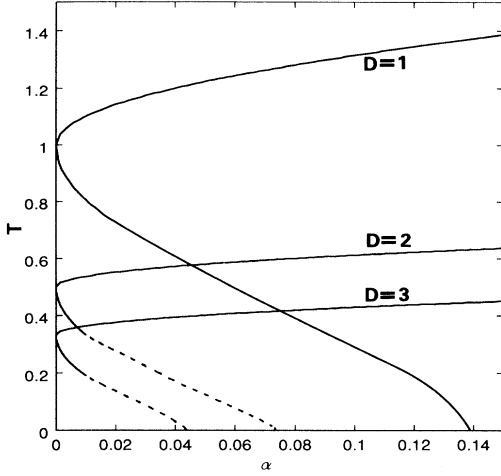
FIG. 5. $T$-$\alpha$ phase diagram of a network in the case of $D = 1, 2,$ and 3.

use of this free energy, first, the case of $\alpha = 0$ was studied. This is the case in which a finite number of patterns are embedded in the network with an infinite number of neurons. We found that the patterns can be retrieved if the temperature $T$ is lower than $1/D$. Next, we calculated the storage capacity in the case of $D = 2$ and $D = 3$ and obtained $\alpha_c = 0.0743$ ($D = 2$) and $\alpha_c = 0.0432$ ($D = 3$) for $T = 0$. In the case of $D = 2$, our model is similar to the Cook's one, but the storage capacity of our model is larger than that of Cook's model. One of the reason is that the synaptic coupling $J_{ij}$ of our model is expressed by a matrix, which is in need of more information. We also calculated the phase diagram and found that the retrieval regime shrank as dimension $D$ increased. It was difficult to calculate the storage capacity for $D \geq 4$, but the dependence of the storage capacity on dimension $D$ is an interesting open problem.

## APPENDIX A: DERIVATION OF EQ. (27)

In this Appendix, we show how to calculate $\langle\langle \rangle\rangle_{\xi_i^\mu}$ in Eq. (27). First, we estimate the exponential in the second line of Eq. (27).

$$\exp\left[\beta \sum_\rho \left\{ m_\rho^{\mu t} \xi_i^\mu x_i^\rho - \frac{1}{2N}({}^t\xi_i^\mu x_i^\rho)^2 \right\}\right]$$

$$= \exp\left[\sum_{j=1}^D X_j \xi_{i(j)}^\mu - \sum_{j=1}^D Y_j (\xi_{i(j)}^\mu)^2 - \sum_{j \neq k}^D Z_{jk} \xi_{i(j)}^\mu \xi_{i(k)}^\mu \right]$$

$$\approx 1 + \sum_{j=1}^D X_j \xi_{i(j)}^\mu + \frac{1}{2}\left(\sum_{j=1}^D X_j \xi_{i(j)}^\mu\right)^2$$

$$- \sum_{j=1}^D Y_j (\xi_{i(j)}^\mu)^2 - \sum_{j \neq k}^D Z_{jk} \xi_{i(j)}^\mu \xi_{i(k)}^\mu, \tag{A1}$$

where

$$X_j \equiv \beta \sum_\rho m_\rho^\mu x_{i(j)}^\rho \sim O\left(\frac{1}{\sqrt{N}}\right), \tag{A2}$$

$$Y_j \equiv \frac{\beta}{2N} \sum_\rho (x_{i(j)}^\rho)^2 \sim O\left(\frac{1}{N}\right), \tag{A3}$$

$$Z_{jk} \equiv \frac{\beta}{2N} \sum_\rho x_{i(j)}^\rho x_{i(k)}^\rho \sim O\left(\frac{1}{N}\right), \tag{A4}$$

and we neglected terms which are smaller than $O(1/N)$ in the Taylor expansion. Using Eq. (25), we easily find the following:

$$\left\langle\left\langle \xi_{i(j)}^\mu \right\rangle\right\rangle_{\xi_i^\mu} = 0 \quad (j = 1, \ldots, D), \tag{A5}$$

$$\left\langle\left\langle \xi_{i(j)}^\mu \xi_{i(k)}^\mu \right\rangle\right\rangle_{\xi_i^\mu} = 0 \quad (j, k = 1, \ldots, D; j \neq k), \tag{A6}$$

$$\left\langle\left\langle (\xi_{i(j)}^\mu)^2 \right\rangle\right\rangle_{\xi_i^\mu} = \frac{1}{D} \quad (j = 1, \ldots, D). \tag{A7}$$

Therefore we can calculate as

$$\prod_{\mu i} \left\langle\left\langle \exp\left[\beta \sum_\rho \left\{ m_\rho^{\mu t} \xi_i^\mu x_i^\rho - \frac{1}{2N}({}^t\xi_i^\mu x_i^\rho)^2 \right\}\right]\right\rangle\right\rangle_{\xi_i^\mu}$$

$$= \prod_{\mu i} \left\langle\left\langle 1 + \sum_{j=1}^D X_j \xi_{i(j)}^\mu + \frac{1}{2}\left(\sum_{j=1}^D X_j \xi_{i(j)}^\mu\right)^2 - \sum_{j=1}^D Y_j (\xi_{i(j)}^\mu)^2 - \sum_{j \neq k}^D Z_{jk} \xi_{i(j)}^\mu \xi_{i(k)}^\mu \right\rangle\right\rangle_{\xi_i^\mu}$$

$$= \prod_{\mu i} \left[ 1 + \frac{1}{2D}\sum_{j=1}^D X_j^2 - \frac{1}{D}\sum_{j=1}^D Y_j \right]$$

$$= \prod_{\mu i} \left[ 1 + \frac{\beta^2}{2D}\sum_{\rho\sigma} m_\rho^\mu m_\sigma^{\mu t} x_i^\rho x_i^\sigma - \frac{\beta n}{2DN} \right]$$

$$
\begin{aligned}
&= \prod_{\mu i} \exp \left[ \ln \left\{ 1 + \frac{\beta^2}{2D} \sum_{\rho\sigma} m_\rho^\mu m_\sigma^{\mu t} x_i^\rho x_i^\sigma - \frac{\beta n}{2DN} \right\} \right] \\
&= \exp \left( -\frac{\beta p n}{2D} \right) \exp \left[ \sum_{\mu i} \frac{\beta^2}{2D} \sum_{\rho\sigma} m_\rho^\mu m_\sigma^{\mu t} x_i^\rho x_i^\sigma \right] ,
\end{aligned}
\tag{A8}
$$

where we used an approximation $\ln(1 + \varepsilon) \approx \varepsilon$ for sufficiently small $\varepsilon$. Equation (A8) is the final form of Eq. (27).

---

[1] D. J. Amit, H. Gutfreund, and H. Sompolinsky, Phys. Rev. Lett. **55**, 1530 (1985).

[2] D. J. Amit, H. Gutfreund, and H. Sompolinsky, Phys. Rev. A **32**, 1007 (1985).

[3] D. J. Amit, H. Gutfreund, and H. Sompolinsky, Ann. Phys. **173**, 30 (1987).

[4] D. J. Amit, *Modeling Brain Function* (Cambridge University Press, Cambridge, 1989).

[5] P. Peretto, *An Introduction to the Modeling of Neural Networks* (Cambridge University Press, Cambridge, 1989).

[6] H. Rieger, J. Phys. A **23**, L1273 (1990).

[7] D. Bollé and J. von Mourik, J. Phys. A **27**, 1151 (1994).

[8] D. Bollé, G. M. Shim, B. Vinck, and V. A. Zagrebnov, J. Stat. Phys. **74**, 565 (1987).

[9] I. Kanter, Phys. Rev. A **37**, 2739 (1988).

[10] D. Bollé, P. Dupont, and J. von Mourik, J. Phys. A **24**, 1065 (1991).

[11] A. J. Noest, Europhys. Lett. **6**, 469 (1988).

[12] A. J. Noest, Phys. Rev. A **38**, 2196 (1988).

[13] J. Cook, J. Phys. A **22**, 2057 (1989).

[14] F. Gerl, K. Bauer, and U. Krey, Z. Phys. B **88**, 339 (1992).

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