

Unlearning of Mixed States in the Hopfield Model – Finite Loading Case –

Haruka Otani¹, Midori Yoshida¹, Shuji Kiyokawa², and Tatsuya Uezu¹ *

¹ *Graduate School of Humanities and Sciences, Nara Women's University, Nara 630, Japan*

² *Faculty of Science, Nara Women's University, Nara 630, Japan*

We study the unlearning of mixed states in the Hopfield model for the finite loading case, that is, $\alpha = \frac{p}{N} \ll 1$, where N and p are the numbers of neurons and embedded patterns, respectively. In the general situation that any number of mixed states that exist in the model is unlearned, we derive the saddle point equations (SPEs) and evolution equations for overlaps by introducing sublattices. We postulate a condition that the solutions are stable in equilibrium, and prove that the static and dynamic stabilities are the same. We also prove that the stable state of the Hopfield model continuously changes and is statically and dynamically stable for sufficiently small unlearning coefficients. For $p = 3$, we perform detailed theoretical and numerical calculations. In the case that a single mixed state is unlearned, we determine phase boundaries using the Hessian matrix and by numerically integrating evolution equations. We performed Markov chain Monte Carlo simulations and find that the simulation results agree with the theoretical ones reasonably well. For general p , when all of the mixed states are unlearned with an equal unlearning coefficient η , we derive the formulae for critical unlearning coefficients at the temperature $T = 0$ below which embedded patterns and mixed states exist and are stable as solutions of the SPEs. We found that there is an unlearning region of (T, η) in which all patterns are retained and all mixed states are deleted, although tuning the parameters in this region is more difficult as p increases since the region shrinks. We numerically confirmed the theoretical results of the p dependences of the critical unlearning coefficients.

1. Introduction

One of the roles of dreaming is considered to be to regulate memories, that is, unnecessary memories are considered to be deleted by dreaming.¹⁾ In this context, we consider strengthening important memories and weakening unnecessary memories. In this work, as a concrete

*uezu@cc.nara-wu.ac.jp

model, we study the Hopfield model.²⁾ As is well known, in the Hopfield model, when p patterns are stored, a combination of several patterns, a mixed state, is also stored. There are several types of stable mixed state.³⁻⁵⁾ We study the unlearning of mixed states. There have been many studies on unlearning in neural networks.⁶⁻¹⁰⁾ These papers treat unlearning the final states starting from noisy inputs, or unlearning by parallel dynamics, or unlearning by using asymmetric synaptic weights, and so forth. In particular, in Refs. 11 and 12, as a candidate of unnecessary memories, a spin glass state is considered.

In this study, we treat symmetric synaptic weights and asynchronous dynamics. In the Hopfield model, spin glass states appear for the extensive loading of patterns. Since we study the finite loading case, mixed states are only spurious states. We investigate the statics and dynamics of the network taking unlearning into account.

We postulate that the stability at the equilibrium is determined by the condition that the para state is stable at high temperatures. We prove that the static stability determined by this assumption agrees with the dynamic stability in general situations. First of all, as the simplest case, we study the case $p = 3$, and perform detailed theoretical and numerical calculations. In particular, for the case that a single mixed state is unlearned, we numerically integrated evolution equations by the Runge-Kutta (RK) method, and performed Markov chain Monte Carlo (MCMC) simulations. The numerical results agree with the theoretical ones quite well. Since we found that the RK method gives the same phase boundaries as those determined by the eigenvalues of the Hessian matrix and MCMC simulations for the unlearning of the single mixed state, we study the unlearning of multiple mixed states by the RK method for general p . Secondly, we study the case that all mixed states are unlearned and derive critical unlearning coefficients at $T = 0$, below which embedded patterns and mixed states exist. We show that we can eliminate all mixed states and retain all patterns if the temperature and the unlearning coefficient take values in some region in (T, η) space. The larger the number of patterns, the smaller the difference between two critical coefficients. In particular, both coefficients are zero at $p = \infty$. This implies that the larger the value of p , the more difficult it is to tune the parameters in the region.

This paper is organized as follows. In sect. 2, we formulate the problem in the general situation that any number of mixed states is unlearned, and describe the saddle point equations (SPEs) of overlaps in the equilibrium and evolution equations for overlaps and sublattice overlaps. In sect. 3, we study in detail the case that one and several mixed states are unlearned for $p = 3$ in detail. In sect. 4, we study the case that all mixed states are unlearned for a general p . Section 5 contains a summary and discussion of the results. In Appendices A and B, we

derive evolution equations for sublattice overlaps and prove the equivalence of the static and dynamic stabilities of the solutions of the SPEs, respectively. In Appendices C and D, in the case that the single mixed state is unlearned, we describe the SPEs and the stability of the solutions of the SPEs, respectively.

2. Formulation

The Hopfield model is a recurrent network of N neurons, in which all neurons interact with each other. The state of the i th neuron is represented by s_i . $s_i = 1$ or $s_i = -1$ corresponds to a firing state or a rest state, respectively. The number of patterns is set to p . Let $\xi^\mu = (\xi_1^\mu, \xi_2^\mu, \dots, \xi_N^\mu)$ be μ th pattern, where $\mu = 1, 2, \dots, p$. We assume that ξ_i^μ takes values of ± 1 independently with the probability $1/2$. Let us denote the configuration of N neurons as $s = \{s_i\}$ and let J_{ij} be the synaptic weight from the j th neuron to the i th neuron. We assume $J_{ij} = J_{ji}$ and $J_{ii} = 0$. The synaptic weight of the Hopfield model, $J_{ij}^{(H)}$, is given by

$$\begin{cases} J_{ij}^{(H)} = \frac{1}{N} \sum_{\mu=1}^p \xi_i^\mu \xi_j^\mu, & (i \neq j), \\ J_{ii}^{(H)} = 0. \end{cases} \quad (1)$$

The input signal to the i th neuron at time t , $h_i(t)$, is given by

$$h_i(t) = \sum_{j(\neq i)} J_{ij} s_j(t). \quad (2)$$

In the deterministic update, the new state of the i th neuron is

$$s_i(t + \Delta t) = \text{sgn}\left(\sum_{j(\neq i)} J_{ij} s_j(t)\right), \quad (3)$$

where $\text{sgn}(x) = 1$ for $x \geq 0$ and -1 for $x < 0$. We take the time increment $\Delta t = \frac{1}{N}$. In this study, we consider the probabilistic update, and the probability that the i th neuron takes the value ± 1 is given by

$$\text{Prob}[s_i(t + \Delta t) = \pm 1] = \frac{1 \pm \tanh[\beta h_i(t)]}{2}. \quad (4)$$

Here, $\beta = \frac{1}{T}$ and T represents the strength of noise and is called the ‘temperature’ in this paper.

2.1 Statics

By adopting an asynchronous update, the stationary state of the network becomes equivalent to the equilibrium state of the canonical ensemble with the following Hamiltonian H :

$$H = - \sum_{i < j} J_{ij} s_i s_j. \quad (5)$$

In the Hopfield model, for $p \ll N$, all patterns are stored when $T < T_c$, where T_c is given by $T_c = 1$. Furthermore, mixed states with three patterns are also stable when $T < T_3$, where $T_3 \approx 0.46$.⁴⁾

In this paper, we treat the case of a general p and consider the unlearning of any number of mixed states with three patterns. A mixed state is denoted by $\xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} = (\xi_1^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}, \xi_2^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}, \dots, \xi_N^{(\boldsymbol{\mu}; \boldsymbol{\gamma})})$ and determined as

$$\xi_i^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} = \text{sgn}(\gamma_1 \xi_i^{\mu_1} + \gamma_2 \xi_i^{\mu_2} + \gamma_3 \xi_i^{\mu_3}), \quad i = 1, \dots, N. \quad (6)$$

Here, $\boldsymbol{\mu} = (\mu_1, \mu_2, \mu_3)$ and μ_1, μ_2 , and μ_3 are all different and take values in $\{1, 2, \dots, p\}$. $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \gamma_3)$ and $|\gamma_i| = 1, (i = 1, 2, 3)$. Without loss of generality, we set $\mu_1 < \mu_2 < \mu_3$. We subtract the mixed states from $J_{ij}^{(H)}$ and treat the following synaptic weight J_{ij} ,

$$J_{ij} = J_{ij}^{(H)} + J_{ij}^{(U)} = \frac{1}{N} \sum_{\mu=1}^p \xi_i^{\mu} \xi_j^{\mu} + \frac{1}{N} \sum_{\mathcal{V}} \zeta(\boldsymbol{\mu}; \boldsymbol{\gamma}) \xi_i^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi_j^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}. \quad (7)$$

Here, $\zeta(\boldsymbol{\mu}; \boldsymbol{\gamma})$ is the unlearning coefficient, that is, when $\zeta(\boldsymbol{\mu}; \boldsymbol{\gamma})$ is negative, $\xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}$ is unlearned. \mathcal{V} is the set of superscripts $(\boldsymbol{\mu}; \boldsymbol{\gamma})$ with which a mixed state $\xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}$ is unlearned, and $\sum_{\mathcal{V}}$ denotes the summation of all mixed states with superscripts in the set \mathcal{V} . Let u be the number of unlearned mixed states, $u \equiv |\mathcal{V}|$. When all patterns ξ^1, \dots, ξ^p change their signs, that is, by the transformation $(\xi^1, \dots, \xi^p) \rightarrow (-\xi^1, \dots, -\xi^p)$, J_{ij} does not change. We call this the inversion symmetry. Because of the inversion symmetry, we may set $\gamma_1 = 1$ without loss of generality. At first sight, Eq. (7) seems to represent the unlearning of the mixed states $\xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}$ only. However, there is a correlation between the mixed state with $\boldsymbol{\mu} = (\mu_1, \mu_2, \mu_3)$ and the pattern ξ^{μ_k} , ($k = 1, 2, 3$), that is,

$$\frac{1}{N} \sum_{i=1}^N \xi_i^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi_i^{\mu_k} = \langle\langle \xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi^{\mu_k} \rangle\rangle = \frac{1}{2} \gamma_k \quad \text{for } k = 1, 2, 3. \quad (8)$$

Here, double brackets $\langle\langle \cdot \rangle\rangle$ denote the average over $\xi^1, \xi^2, \dots, \xi^p$, where ξ^{μ} takes values of ± 1 with probability 1/2. In this paper, when a quantity is averaged and expressed by double brackets, we omit its subscript, which is i in the present case, if it does not cause any confusion. The first equality follows from the self-averaging because $N \gg 2^p$. Therefore, a part of ξ^{μ_k} is subtracted from J_{ij} for $k = 1, 2, 3$. Nevertheless, we still call the present procedure the ‘unlearning’ of the mixed states for simplicity. For later use, we write correlations between two mixed states $\xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}$ and $\xi^{(\boldsymbol{\mu}'; \boldsymbol{\gamma}')}$ where $\boldsymbol{\mu}' = (\mu'_1, \mu'_2, \mu'_3)$ and $\boldsymbol{\gamma}' = (\gamma'_1, \gamma'_2, \gamma'_3)$.

$$\frac{1}{N} \sum_i \xi_i^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi_i^{(\boldsymbol{\mu}'; \boldsymbol{\gamma}')} = \langle\langle \xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi^{(\boldsymbol{\mu}'; \boldsymbol{\gamma}')} \rangle\rangle \quad (9)$$

$$= \begin{cases} 0 & \{\mu_1, \mu_2, \mu_3\} \cap \{\mu'_1, \mu'_2, \mu'_3\} = \phi, \text{ none of the elements of } \boldsymbol{\mu} \\ & \text{and } \boldsymbol{\mu}' \text{ agree.} \\ \frac{1}{4}\gamma_i\gamma'_i, & \mu_i = \mu'_i, \text{ one of the elements of } \boldsymbol{\mu} \text{ and } \boldsymbol{\mu}' \text{ agrees.} \\ \frac{1}{4}(\gamma_i\gamma'_i + \gamma_j\gamma'_j), & \mu_i = \mu'_i, \mu_j = \mu'_j, \text{ two of the elements of } \boldsymbol{\mu} \text{ and } \boldsymbol{\mu}' \text{ agree.} \\ 1 & \boldsymbol{\mu} = \boldsymbol{\mu}', \boldsymbol{\gamma} = \boldsymbol{\gamma}' \\ 0 & \boldsymbol{\mu} = \boldsymbol{\mu}', \boldsymbol{\gamma} \neq \boldsymbol{\gamma}' \end{cases} \quad (10)$$

2.2 SPEs

In this subsection, we formulate the general case that any number of mixed states is unlearned. For simplicity, we denote mixed states by $\boldsymbol{\xi}^{p+1}, \dots, \boldsymbol{\xi}^v$, where $v = p + u$, and rewrite J_{ij} [Eq. (7)] as

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^v \zeta_{\mu} \xi_i^{\mu} \xi_j^{\mu}, \quad (11)$$

where $\zeta_1 = \dots = \zeta_p = 1$. The order parameters of the system are the overlaps m^{μ} between \mathbf{s} and $\boldsymbol{\xi}$ s,

$$m^{\mu}(t) = \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu} s_i(t), \quad (\mu = 1, \dots, v). \quad (12)$$

By the standard recipe, the free energy per neuron f and the SPEs are derived as

$$f = \sum_{\mu=1}^v \zeta_{\mu} \frac{(m^{\mu})^2}{2} - \frac{1}{\beta} \langle \ln[2 \cosh \beta (\sum_{\mu=1}^v \zeta_{\mu} m^{\mu} \xi^{\mu})] \rangle, \quad (13)$$

$$m^{\mu} = \langle \langle \xi^{\mu} \tanh[\beta (\sum_{\nu=1}^v \zeta_{\nu} m^{\nu} \xi^{\nu})] \rangle \rangle, \quad \mu = 1, \dots, v. \quad (14)$$

2.3 Stability of solutions of the SPEs

In this subsection, we study the stability of the solutions of the SPEs. The stability of a solution is determined by the Hessian matrix $\Lambda = \{\Lambda_{\mu\nu}\}$, where $\Lambda_{\mu\nu} = \frac{\partial^2 f}{\partial m^{\mu} \partial m^{\nu}}$, which is evaluated at the solution. To find a condition under which the solution is stable, we study the behavior of f at high temperatures. For $\beta \ll 1$, f is expressed as

$$f \simeq \sum_{\mu=1}^v \zeta_{\mu} \frac{(m^{\mu})^2}{2} + \text{const.} \quad (15)$$

Since the para (P) state in which the m^{μ} are 0 for all μ should be stable at high temperatures, f should be minimal with respect to m^{μ} for positive ζ_{μ} and be maximal with respect to m^{μ} for negative ζ_{μ} . Let λ be the eigenvalues of Λ , and let us define λ as $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \lambda_4 \dots \geq \lambda_v$. For example, let us consider the case that $p = 3$ and only the mixed state $\boldsymbol{\xi}^4 \equiv \boldsymbol{\xi}^{(1,2,3;1,1,1)}$ is

unlearned. Defining $\zeta_1 = \zeta_2 = \zeta_3 = 1$ and $\zeta_4 = -\eta$, the condition under which the solution is stable is given by $\lambda_1 > 0, \lambda_2 > 0, \lambda_3 > 0, \lambda_4 < 0$ for $\eta > 0$ and $\lambda_1 > 0, \lambda_2 > 0, \lambda_3 > 0, \lambda_4 > 0$ for $\eta < 0$. Thus, in general situations, a solution is stable when the number of positive (negative) eigenvalues is the same as the number of positive (negative) ζ_μ . We prove that the thus defined static stability agrees with the dynamic stability in Appendix B.

In the next subsection, we formulate the dynamics of the network.

2.4 Dynamics

Let $p_t(\mathbf{s})$ be the probability that the configuration of neurons is $\mathbf{s} = (s_1, s_2, \dots, s_k, \dots, s_N)$ at time t and let $w_k(\mathbf{s})$ be the transition probability from \mathbf{s} to $F_k \mathbf{s} \equiv (s_1, s_2, \dots, -s_k, \dots, s_N)$ per unit of time. F_k is the flip operator of the k th neuron. We adopt the following function as $w_k(\mathbf{s})$:

$$w_k(\mathbf{s}) = \frac{1}{2}[1 - s_k \tanh(\beta h_k)], \quad (16)$$

$$h_k = \sum_{j(\neq k)} J_{kj} s_j. \quad (17)$$

Then, the time evolution of $p_t(\mathbf{s})$ is described by the following master equation:

$$\frac{\partial}{\partial t} p_t(\mathbf{s}) = \sum_{k=1}^N [w_k(F_k \mathbf{s}) p_t(F_k \mathbf{s}) - w_k(\mathbf{s}) p_t(\mathbf{s})]. \quad (18)$$

It is shown that the stationary state of Eq. (18) is the equilibrium state of the canonical ensemble with the Hamiltonian given by Eq. (5).

Now, let us introduce sublattices and sublattice overlaps. The number of possible configurations of $\{\xi_i^\mu\} \equiv (\xi_i^1, \xi_i^2, \dots, \xi_i^p)$ with i fixed is 2^p . For each configuration, we introduce a sublattice Λ_l of neurons in such a way that $\Lambda_{l+2^{p-1}} = \Lambda_l$ holds for $l = 1, \dots, 2^{p-1}$. The number of elements in Λ_l , $|\Lambda_l|$, is $|\Lambda_l| = \frac{N}{2^p}$ ($l = 1, 2, \dots, 2^p$). The sublattice overlap \mathcal{M}_l is defined as

$$\mathcal{M}_l = \frac{1}{|\Lambda_l|} \sum_{i \in \Lambda_l} s_i = \frac{2^p}{N} \sum_{i \in \Lambda_l} s_i, \quad l = 1, \dots, 2^p.$$

Using $\{\mathcal{M}_l\}$, the overlap m^μ is expressed as

$$m^\mu = \frac{1}{N} \sum_{i=1}^N \xi_i^\mu s_i = \frac{1}{N} \sum_{l=1}^{2^p} \sum_{i \in \Lambda_l} \xi_i^\mu s_i = \frac{1}{N} \sum_{l=1}^{2^p} \xi^{\mu,l} \sum_{i \in \Lambda_l} s_i = \frac{1}{2^p} \sum_{l=1}^{2^p} \xi^{\mu,l} \mathcal{M}_l, \quad \mu = 1, \dots, p \quad (19)$$

$\xi^{\mu,l}$ is the value of ξ_i^μ for i in Λ_l . Now, let us study the evolution equations for $\{\mathcal{M}_l\}$. We define $\mathbf{M} = (\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_{2^p})$ and $d\mathbf{M} = \prod_{l=1}^{2^p} d\mathcal{M}_l$. Let $p_t(\mathbf{M})$ be the probability density that the sublattice overlap \mathcal{M}_l takes a value in $(\mathcal{M}_l, \mathcal{M}_l + d\mathcal{M}_l)$ for $l = 1, \dots, 2^p$ at time t ,

$$p_t(\mathbf{M}) = \text{Tr}_{\mathbf{s}} p_t(\mathbf{s}) \prod_{l'=1}^{2^p} \delta(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i). \quad (20)$$

The evolution equation for \mathcal{M}_l is

$$\frac{d\mathcal{M}_l}{dt} = -\mathcal{M}_l + \tanh[\beta h^l(\mathcal{M})], \quad l = 1, \dots, 2^p. \quad (21)$$

Here, $h^l(\mathcal{M})$ is expressed as

$$h^l(\mathcal{M}) = \sum_{\mu=1}^v \zeta_{\mu} \Xi_{\mu,l} \frac{1}{2^p} (\Xi \mathcal{M})_{\mu},$$

where Ξ is a $v \times 2^p$ matrix whose (μ, l) element is $\Xi_{\mu,l} = \xi^{\mu,l}$. See Appendix A for the derivation.

For $l \leq 2^{p-1}$, $h^{l+2^{p-1}}(\mathcal{M}) = -h^l(\mathcal{M})$ holds because $\Xi_{\mu,l+2^{p-1}} = -\Xi_{\mu,l}$. Thus, we have

$$\frac{d\mathcal{M}_{l+2^{p-1}}}{dt} = -\mathcal{M}_{l+2^{p-1}} - \tanh[\beta h^l(\mathcal{M})], \quad l = 1, \dots, 2^{p-1}, \quad (22)$$

$$\frac{d\mathcal{M}_l}{dt} = -\mathcal{M}_l + \tanh[\beta h^l(\mathcal{M})], \quad l = 1, \dots, 2^{p-1}. \quad (23)$$

From these equations, we obtain

$$\frac{d(\mathcal{M}_{l+2^{p-1}} + \mathcal{M}_l)}{dt} = -(\mathcal{M}_{l+2^{p-1}} + \mathcal{M}_l).$$

Its solution is

$$\mathcal{M}_l(t) + \mathcal{M}_{l+2^{p-1}}(t) = e^{-t}[\mathcal{M}_l(0) + \mathcal{M}_{l+2^{p-1}}(0)].$$

As $t \rightarrow \infty$, we have

$$\mathcal{M}_{l+2^{p-1}} = -\mathcal{M}_l \text{ for } l \leq 2^{p-1}.$$

Therefore, we assume $\mathcal{M}_{l+2^{p-1}} = -\mathcal{M}_l$ for $l \leq 2^{p-1}$ in equilibrium. Thus, m_{μ} is expressed as

$$m^{\mu} = \frac{1}{2^{p-1}} \sum_{l=1}^{2^{p-1}} \xi^{\mu,l} \mathcal{M}_l, \quad \mu = 1, \dots, v. \quad (24)$$

Evolution equations for m^{μ} are also derived. Let us define the probability density $p_t(\mathbf{m})$ of overlaps $\mathbf{m} = (m^1, m^2, \dots, m^v)$ at time t as

$$p_t(\mathbf{m}) = \text{Tr}_S p_t(s) \prod_{\mu=1}^v \delta(m^{\mu} - \frac{1}{N} \sum_j \xi_j^{\mu} s_j). \quad (25)$$

By using the same recipe as that used to derive Eq. (21), we obtain

$$\frac{d}{dt} m^{\mu} = -m^{\mu} + \langle\langle \xi^{\mu} \tanh[\beta (\sum_{\nu=1}^v \zeta_{\nu} m^{\nu} \xi^{\nu})] \rangle\rangle, \quad \mu = 1, \dots, v. \quad (26)$$

Any stationary state of the evolution equations given by Eq. (26) satisfies the SPEs (14).

2.5 Relationship between static and dynamic stabilities of equilibrium solution

As is proved in Appendix B, the static and dynamic stabilities of any stationary state are the same, and any stable state of the Hopfield model ($\eta = 0$) continuously changes and is stat-

ically and dynamically stable for sufficiently small unlearning coefficients. Furthermore, it is proved in Appendix B that the stationary solution becomes unstable statically and dynamically at the same time in generic situations that all $\zeta_\mu (\mu = 1, \dots, \nu)$ are different. Hereafter, we assume $\zeta_1 = \dots = \zeta_p = 1$ and $\zeta_{p+1} = \dots = \zeta_\nu = -\eta$.

3. Case of $p = 3$

In this section, we study the case of $p = 3$. Sublattices are given as

$$\begin{aligned}
 \Lambda_1 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (1, 1, 1, 1, 1, 1, -1)], \\
 \Lambda_2 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (1, 1, -1, 1, 1, -1, 1)], \\
 \Lambda_3 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (1, -1, 1, 1, -1, 1, 1)], \\
 \Lambda_4 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (1, -1, -1, -1, 1, 1, 1)], \\
 \Lambda_5 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (-1, -1, -1, -1, -1, -1, 1)], \\
 \Lambda_6 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (-1, -1, 1, -1, -1, 1, -1)], \\
 \Lambda_7 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (-1, 1, -1, -1, 1, -1, -1)], \\
 \Lambda_8 &= [i(\xi_i^1, \xi_i^2, \xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6, \xi_i^7) = (-1, 1, 1, 1, -1, -1, -1)],
 \end{aligned} \tag{27}$$

where $\xi_i^4, \xi_i^5, \xi_i^6$, and ξ_i^7 are defined as

$$\xi_i^4 \equiv \xi_i^{(1,2,3:1,1,1)} = \text{sgn}(\xi_i^1 + \xi_i^2 + \xi_i^3), \tag{28}$$

$$\xi_i^5 \equiv \xi_i^{(1,2,3:1,1,-1)} = \text{sgn}(\xi_i^1 + \xi_i^2 - \xi_i^3), \tag{29}$$

$$\xi_i^6 \equiv \xi_i^{(1,2,3:1,-1,1)} = \text{sgn}(\xi_i^1 - \xi_i^2 + \xi_i^3), \tag{30}$$

$$\xi_i^7 \equiv \xi_i^{(1,2,3:1,-1,-1)} = \text{sgn}(\xi_i^1 - \xi_i^2 - \xi_i^3). \tag{31}$$

Correspondingly, m^4, m^5, m^6 , and m^7 are defined as

$$m^4 = \frac{1}{N} \sum_{i=1}^N \xi_i^{(1,2,3:1,1,1)} s_i(t), \tag{32}$$

$$m^5 = \frac{1}{N} \sum_{i=1}^N \xi_i^{(1,2,3:1,1,-1)} s_i(t), \tag{33}$$

$$m^6 = \frac{1}{N} \sum_{i=1}^N \xi_i^{(1,2,3:1,-1,1)} s_i(t), \tag{34}$$

$$m^7 = \frac{1}{N} \sum_{i=1}^N \xi_i^{(1,2,3:1,-1,-1)} s_i(t). \tag{35}$$

Let Ξ and X be the following 7×8 and 7×4 matrices, respectively, whose (μ, l) component is $\xi^{\mu,l}$.

$$\Xi = (X, -X), \quad (36)$$

$$X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & -1 \\ 1 & 1 & -1 & 1 \\ 1 & -1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}. \quad (37)$$

Equation (24) in equilibrium is written as

$$\begin{pmatrix} m^1 & m^2 & m^3 & m^4 & m^5 & m^6 & m^7 \end{pmatrix}^T = \frac{1}{4} X \begin{pmatrix} \mathcal{M}_1 & \mathcal{M}_2 & \mathcal{M}_3 & \mathcal{M}_4 \end{pmatrix}^T,$$

where T implies the transpose. Let X_4 be the 4×4 matrix whose (μ, l) element is $\xi^{\mu,l}$ for $\mu = 1, \dots, 4$, $l = 1, \dots, 4$. Since $|X_4| = 8$, X has an inverse, that is, m^1, \dots, m^4 and $\mathcal{M}_1, \dots, \mathcal{M}_4$ have a one-to-one correspondence. Therefore, in the equilibrium state, m^5, m^6 , and m^7 are expressed in terms of m^1, \dots, m^4 . X_4^{-1} is

$$X_4^{-1} = \frac{1}{2} \begin{pmatrix} 1 & 1 & 1 & -1 \\ 0 & 0 & -1 & 1 \\ 0 & -1 & 0 & 1 \\ 1 & 0 & 0 & -1 \end{pmatrix}. \quad (38)$$

Thus, we have the following relations in equilibrium:

$$\begin{pmatrix} \mathcal{M}_1 \\ \mathcal{M}_2 \\ \mathcal{M}_3 \\ \mathcal{M}_4 \end{pmatrix} = 4X_4^{-1} \begin{pmatrix} m^1 \\ m^2 \\ m^3 \\ m^4 \end{pmatrix} = 2 \begin{pmatrix} m^1 + m^2 + m^3 - m^4 \\ -m^3 + m^4 \\ -m^2 + m^4 \\ m^1 - m^4 \end{pmatrix}, \quad (39)$$

$$\begin{pmatrix} m^5 \\ m^6 \\ m^7 \end{pmatrix} = \begin{pmatrix} m^1 + m^2 - m^4 \\ m^1 + m^3 - m^4 \\ -m^2 - m^3 + m^4 \end{pmatrix}. \quad (40)$$

From Eq. (19), it follows that the dynamic stability of $\{m^\mu\}$ and that of $\{\mathcal{M}_l\}$ are the same. Also, from Eq. (39), it follows that the static stability of $\{m^\mu\}_{\mu=1,\dots,4}$ and that of $\{\mathcal{M}_l\}_{l=1,\dots,4}$ are the same.

3.1 Case that one mixed state is unlearned

In this subsection, we study the case that only $\xi^4 = \xi^{(1,2,3;1,1,1)}$ is unlearned. The SPEs are

$$m^\mu = \langle \langle \xi^\mu \tanh(\beta W_i) \rangle \rangle, \mu = 1, \dots, 4, \quad (41)$$

$$W_i = \sum_{\mu=1}^3 \xi_i^\mu m^\mu - \eta \xi_i^4 m^4. \quad (42)$$

Here, we put $\zeta_4 = -\eta$. Let us describe the solutions of the SPEs. Let (m^1, m^2, m^3, m^4) be a solution of the SPEs (41). Then, $(m^{1'}, m^{2'}, m^{3'}, m^4)$, in which $m^{1'}$, $m^{2'}$, and $m^{3'}$ are a permutation of m^1 , m^2 , and m^3 , is also a solution of the SPEs (41). We call this the permutation symmetry. Furthermore, $-\mathbf{m}$ is also a solution of the SPEs (41). This is the inversion symmetry. Because of the inversion and permutation symmetries, we assume $m^1 > 0$ and $m^1 \geq m^2 \geq m^3$ without loss of generality. The following types of solution of the SPEs exist.

- (1) S₂ state. $m^1 > m^2 = m^3$.

The Hopfield attractor H characterized by $\mathbf{m} = (m^1, 0, 0, m^4)$ does not exist when the unlearning term exists. See Appendix C. Since $m^2 = m^3 \simeq 0$ in the S₂ state for $|\eta| \ll 1$, S₂ is regarded as a variant of the Hopfield attractor.

- (2) M₄ state. $m^1 = m^2 = m^3$.

This corresponds to the mixed state M₄^(H) in the Hopfield model characterized by $m^1 = m^2 = m^3$.

- (3) M₅ state. $m^1 = m^2 \simeq -m^3$.

The mixed state M₅^(H) in the Hopfield model characterized by $m^1 = m^2 = -m^3$ does not exist when the unlearning term exists. See Appendix C. The M₅ state is considered to be a variant of M₅^(H).

- (4) S₃ state. m^1, m^2 , and m^3 are all different.

- (5) P state. $\mathbf{m} = (0, 0, 0, 0)$.

In the Hopfield model, there are other mixed states M₆^(H) and M₇^(H) characterized by $m^1 = -m^2 = m^3$ and $m^1 = -m^2 = -m^3$, respectively. In the present model, t other mixed states exist: M₆ characterized by $m^1 = m^3 \simeq -m^2$ corresponding to M₆^(H) and M₇ characterized by $m^1 \simeq -m^2 = -m^3$ corresponding to M₇^(H). Because of the permutation and inversion symmetries, these mixed states M₆ and M₇ have the same static and dynamic stabilities as the M₅ state. Thus, we only have to study the M₄ and M₅ states.

Now, let us show some numerical results. We performed MCMC simulations. The transition probability from $s_k \rightarrow -s_k$ is given by Eq. (16).

3.1.1 Statics

As initial conditions, we chose a Hopfield-attractor-like configuration $\mathbf{m} = (m^1, 0, 0, m^4)$ and a mixed-state-like configuration $\mathbf{m} = (m^1, m^1, m^1, m^4)$. It turned out that N should be of the order 10^5 so that the theoretical and simulation results agree. We took several values of N and several hundred Monte Carlo sweeps (MCSs), where one MCS corresponds to one update of N neurons on average. The overlaps were calculated by averaging the overlaps in the time period from 450 to 500 MCSs. The typical number of samples is 10. In Fig. 1, we display the

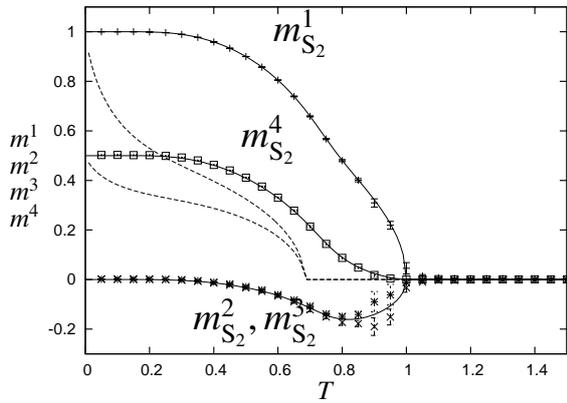


Fig. 1. Temperature dependences of overlaps. $\eta = 0.5$. Curves: theory. Solid curve: S_2 , dashed curve: M_4 . Symbols: Monte Carlo simulations. $N = 10^5$. Averages are taken from 10 samples. Vertical lines denote the error bars. $+$: m^1 , \times : m^2 , $*$: m^3 , \square : m^4 . Initial configuration: $\mathbf{m} = (0.5, 0.5, 0.5, 1)$.

temperature dependences of overlaps for $\eta = 0.5$. In both the Hopfield-attractor-like initial configuration (for which results are not shown) and the mixed-state-like initial configuration, the S_2 state appears for $0 < T < 0.85$ and the S_3 state appears for $0.85 < T \leq 1$. Theoretically, we could not find the S_3 state. Thus, it seems that the theoretical and simulation results do not agree, at least for $0.85 < T \leq 1$. We will discuss this in the next subsection. In Figs. 2(a) and 2(b), we display the temperature dependences of overlaps for $\eta = -0.5$. Starting from a Hopfield-attractor-like initial configuration, the S_2 state appears for $0 < T < 0.57$ and the mixed state M_4 appears for $0.58 < T < 1.42$, as shown in Fig. 2(a). On the other hand, starting from a mixed-state-like initial configuration, the mixed state M_4 appears for $0 < T < 1.42$. See Fig. 2(b). The S_2 state and the unstable S_2 state, which is not drawn in the figure, are annihilated at $T \simeq 0.58$. For $0 < T < 0.57$, the Hopfield attractor and the mixed state M_4 coexist.

The agreement between the numerical and theoretical results is quite good. In the next subsection, we study whether the S_3 state exists for $\eta = 0.5$ at approximately $T = 1$.

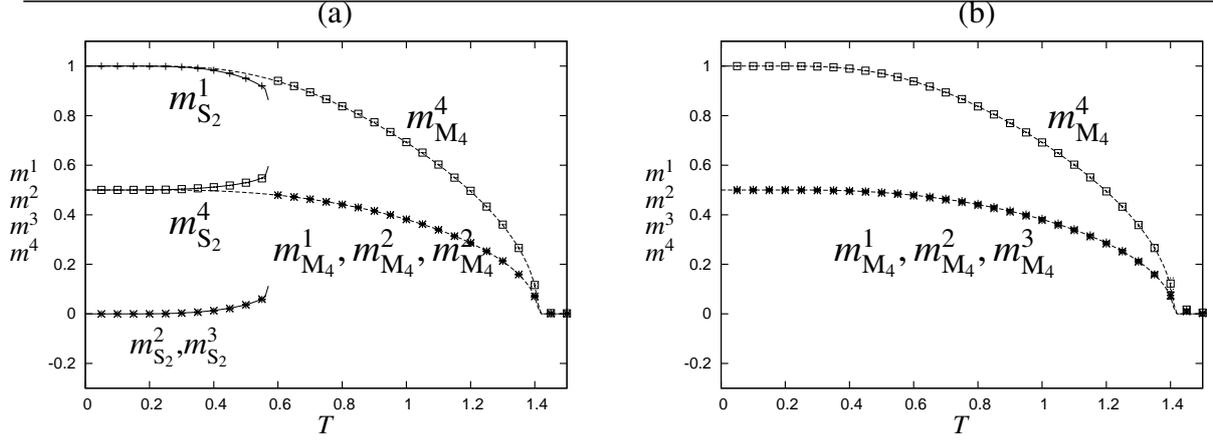


Fig. 2. Temperature dependences of overlaps. $\eta = -0.5$. Curves: theory. Solid curve: S_2 , dashed curve: M_4 . Symbols: Monte Carlo simulations. $N = 10^5$. Averages are taken from 10 samples. Vertical lines denote the error bars. +: m^1 , x: m^2 , *: m^3 , □: m^4 . (a) Initial configuration: $\mathbf{m} = (1, 0, 0, 0.5)$. (b) Initial configuration: $\mathbf{m} = (0.5, 0.5, 0.5, 1)$.

3.1.2 Search for S_3 state near critical temperature, $T \simeq T_c$

We performed MCMC simulations for larger values of the system size N and calculated the difference $m^2 - m^3$ at $T = 0.9$ and 0.95 . In Fig. 3, for $T = 0.9$, we display $\log_{10} N$ vs $\log_{10}(m^2 - m^3)$ with error bars for $N = 10^5, N = 10^6, N = 3 \times 10^6$, and $N = 10^7$ with numbers of samples of 100, 50, 50, and 20, respectively. We estimated $\log_{10}(m^2 - m^3) \sim$

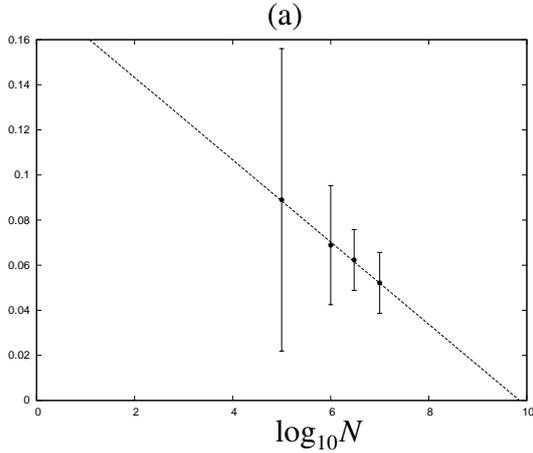


Fig. 3. $\log_{10} N$ vs $\log_{10}(m^2 - m^3)$. $T = 0.9$.

$-0.009926 \cdot \log_{10} N + 0.127675$ at $T = 0.9$ and $\log_{10}(m^2 - m^3) \sim -0.047090 \cdot \log_{10} N + 0.328158$ at $T = 0.95$. From these results, we note that the difference between m^2 and m^3 decreases as N increases. That is, the numerical results of observing the S_3 state is a finite size effect and it is concluded that the S_2 state appears near T_c for $N = \infty$ as the theory predicts.

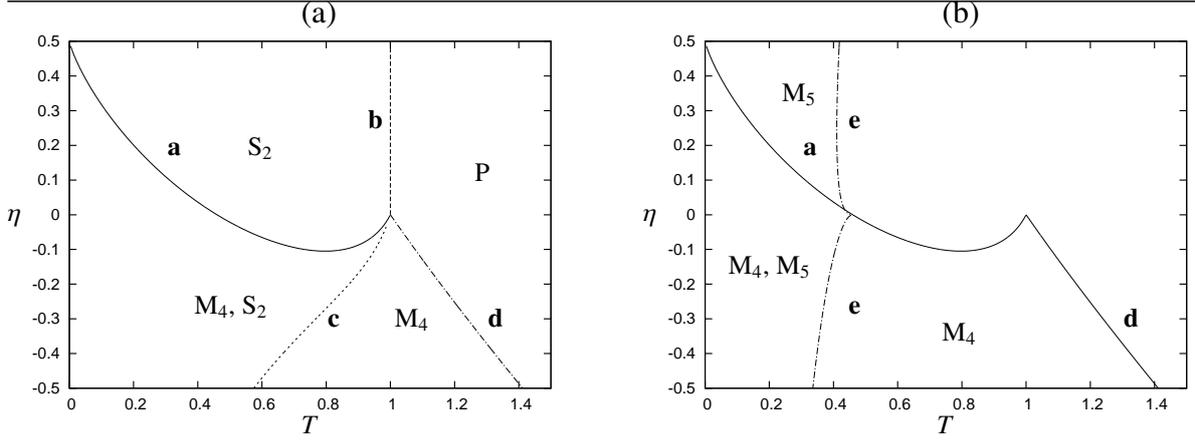


Fig. 4. Phase diagram in (T, η) space determined by using Hessian matrix. (a) Stable regions for P, S_2 , and M_4 . (b) Stable regions for M_4 and M_5 .

3.1.3 Phase diagram

Now, let us study the phase diagram in (T, η) space. The stable regions for solutions are determined by using the Hessian matrix. That is, the boundaries **a**, \dots , **e** in Fig. 4 are determined as follows.

a : Boundary between the S_2 state and the coexistent region of the S_2 and M_4 states. From Fig. 1, we note that at low temperatures, the M_4 state exists but is unstable. The boundary is determined by an eigenvalue of the Hessian matrix for the M_4 state. See Eq. (D·11) in Appendix D.

$$\lambda = \Lambda_{11} - \Lambda_{12} = 0, \text{ i.e., } \frac{\beta}{\cosh^2[\beta(m^1 - \eta m^4)]} = 1.$$

b : Boundary between the P and S_2 states for $\eta > 0$, which is determined by an eigenvalue of the Hessian matrix for the P state. See Eq. (D·7) in Appendix D.

$$\lambda_1 = \lambda_2 = \Lambda_{11} = 1 - \beta = 0, \text{ i.e., } T = 1.$$

c : Boundary between the M_4 state and the coexistent region of the S_2 and M_4 states. This is determined by an eigenvalue of the Hessian matrix for the S_2 state. See Eq. (D·17) in Appendix D. From Fig. 2(a), we note that at this boundary, stable and unstable S_2 states merge and are annihilated.

$$\lambda_4 = 0, \text{ i.e., } -\Lambda_{23}\Lambda_{11}\Lambda_{44} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{12}\Lambda_{24}\Lambda_{14} + 2\Lambda_{11}\Lambda_{24}^2 + 2\Lambda_{44}\Lambda_{12}^2 + \Lambda_{14}^2\Lambda_{23} + \Lambda_{11}\Lambda_{14}^2 = 0.$$

d : Boundary between the P and M_4 states for $\eta < 0$, which is determined by the eigenvalue λ_- of the Hessian matrix for the P state. See Eq. (D·8) in Appendix D.

$$\lambda_- = \frac{1}{2} \left(\Lambda_{11} + \Lambda_{44} - \sqrt{(\Lambda_{11} - \Lambda_{44})^2 + 12\Lambda_{14}^2} \right) = 0, \text{ i.e., } \eta = \frac{4T(1-T)}{4T-1}.$$

e : Boundary of the stable M_5 state. The boundary is determined by

$$-\Lambda_{11}\Lambda_{44}\Lambda_{12} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{34}\Lambda_{13}\Lambda_{14} + 2\Lambda_{11}\Lambda_{14}^2 + 2\Lambda_{44}\Lambda_{13}^2 + \Lambda_{12}\Lambda_{34}^2 + \Lambda_{11}\Lambda_{34}^2 = 0.$$

See (D-14) in Appendix D.

In Fig. 4, the characters denoted in the regions that are surrounded by the curves represent the stable phases. For $\eta > 0$, by unlearning, the region where the mixed state M_4 is stable decreases in size as η increases, and disappears at $\eta = 0.5$, which is easily proved. On the other hand, the S_2 state exists up to $T = 1$ for $\eta > 0$. Since the S_2 state is regarded as a variant of the Hopfield attractor, it seems that the region where the pattern is stable does not decrease. For $\eta < 0$, as η decreases, the stable region of M_4 increases, but the region where the S_2 state is stable decreases slowly compared with the decrease in the stable region of M_4 for $\eta > 0$ and disappears at about $\eta \sim -1.73$. Figure 4(b) shows the phase diagram only for M_4 and M_5 . The stable region of M_5 is almost unchanged for $\eta > 0$ and decreases slowly as $|\eta|$ increases for $\eta < 0$. That is, the unlearning of ξ^4 does not cause any significant change in size in the stable region of the mixed state M_5 . This is considered to be because $\langle\langle \xi^4 \xi^5 \rangle\rangle = 0$.

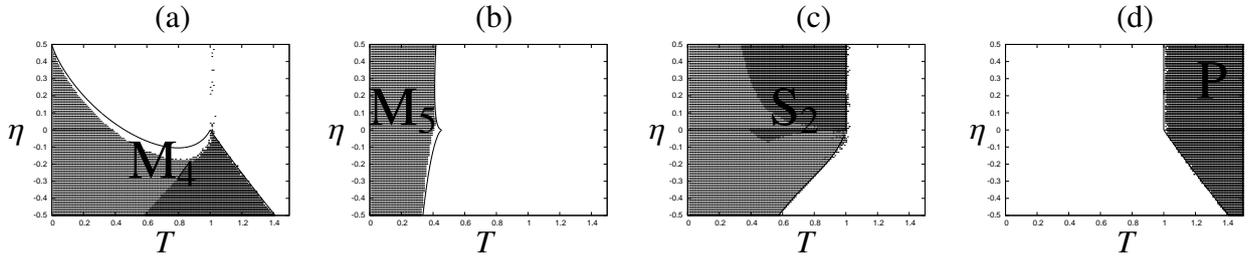


Fig. 5. Stable regions of solutions in (T, η) space. Curves: boundaries obtained by theory using the Hessian matrix, dots: stable regions obtained by MCMC simulation. $N = 10^5$.

To confirm the stable regions for solutions numerically, we performed MCMC simulations and numerically integrated Eq. (26) by the RK method. We iterated the RK routine 10000 times with a time increment of 0.01. This implies that we integrate the evolution equations up to $t = 100$, which is considered to be sufficient for trajectories to converge (as seen later in Fig. 7). On the other hand, we performed MCMC simulations for 10^6 MCSs and used the following criteria to decide the convergent state. First of all, we selected the maximum value of $|m^\mu|$ among $|m^1|, |m^2|$, and $|m^3|$, and if $m^\mu < 0$, we changed the signs of all overlaps. Then, we renumbered the overlaps as $m^1 > m^2 > m^3$. We started with both the Hopfield-attractor-like and mixed-state- M_4 -like initial conditions. We successively imposed the following conditions to decide the convergent states.

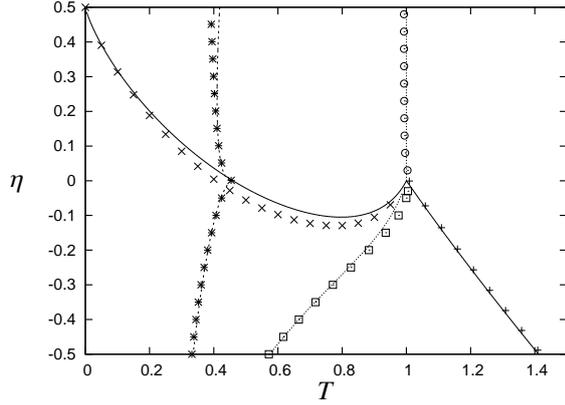


Fig. 6. Phase boundaries obtained by theory using the Hessian matrix (curves) and by the RK method (symbols).

- (1) If $m^1 < 0.04$, the state is regarded as P,
- (2) else if $m^1 - m^2 < 0.08$ and $m^1 - m^3 < 0.08$, the state is regarded as M_4 ,
- (3) else if $|m^2| < 0.005$ and $|m^3| < 0.005$, the state is regarded as H,
- (4) else if $|m^2 - m^3| < 0.01$, the state is regarded as S_2 ,
- (5) else, the state is regarded as S_3 .

The results obtained for both initial conditions are used to decide the convergent states. On the other hand, to determine the stable region of M_5 , we started the M_5 -like initial condition and added the following criterion after (3),

- (3') else if $|m^1 - m^2| < 0.25$ and $|m^1 - m^3| < 0.25$, the state is regarded as M_4 ,

In Fig. 5, we display the theoretical results obtained by using the Hessian matrix and the numerical results obtained by MCMC simulations. In Fig. 6, we draw the phase boundaries obtained using the Hessian matrix and by the RK method. On the boundary of the M_4 state for $T < 1$, the results obtained by MCMC simulations and using the Hessian matrix in Fig. 5(a) do not agree very well. The results obtained by the RK method and using the Hessian matrix in Fig. 6 also do not agree very well. We consider that the reason for this is that the criteria for determining the boundaries of stable regions of solutions by MCMC simulations have ambiguities, that is, we chose conditions such as 0.04 in criterion (1) heuristically. Except for at this boundary, the theoretical and numerical results agree reasonably well.

3.1.4 Dynamics

We performed MCMC simulations for $N = 10^5$ by taking Hopfield-attractor-like and mixed-state-like initial configurations. We compare simulation results obtained by MCMC

simulations with theoretical ones obtained by the RK method. In Fig. 7(a), we display the time series of overlaps for $\eta = 0.5$ and $T = 0.2$. Only the S_2 state exists in this region. Under both initial conditions, all trajectories tend to the S_2 state. In Fig. 7(b), we display the time series of overlaps for $\eta = -0.5$ and $T = 0.2$. In this region, S_2 and M_4 coexist. Depending on the initial conditions, trajectories tend to the S_2 state and the mixed state M_4 . The theoretical and simulation results agree quite well.

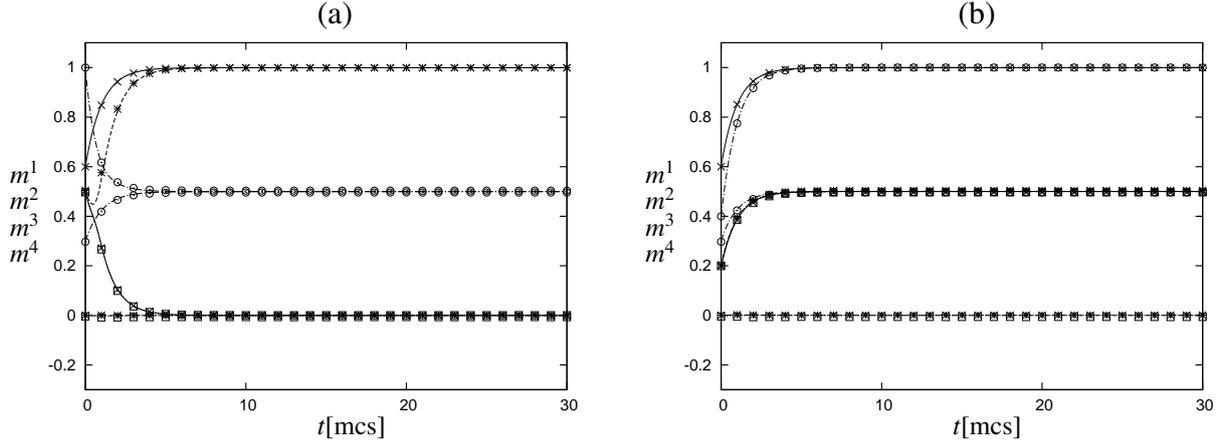


Fig. 7. Curves: theory (RK). Solid curve: m^1 , dashed curve: m^2 , dotted curve: m^3 , dashed-dotted curve: m^4 . Symbols: Monte Carlo simulations. $N = 10^5$. \times : m^1 , $*$: m^2 , \square : m^3 , \circ : m^4 . (a) $\eta = 0.5$, $T = 0.2$. Initial conditions are $\mathbf{m} = (0.6, 0, 0, 0.3)$ and $\mathbf{m} = (0.49, 0.5, 0.49, 1)$. (b) $\eta = -0.5$, $T = 0.2$. Initial conditions are $\mathbf{m} = (0.6, 0, 0, 0.3)$ and $\mathbf{m} = (0.196, 0.204, 0.199, 0.4)$.

The case that another mixed state is unlearned is reduced to the unlearning of ξ^4 as follows. Let us consider the unlearning of $\xi^{(1,2,3;\gamma_1,\gamma_2,\gamma_3)}$. Defining quantities as $\xi^{\mu'} = \gamma_\mu \xi^\mu$, $m^{\mu'} = \gamma_\mu m^\mu$ ($\mu = 1, 2, 3$), and $\xi_i^{4'} = \xi_i^{(1,2,3;\gamma_1,\gamma_2,\gamma_3)} = \text{sgn}(\xi_i^{1'} + \xi_i^{2'} + \xi_i^{3'})$, $m^{4'} = m^{(1,2,3;\gamma_1,\gamma_2,\gamma_3)}$, the SPEs become

$$m^{\mu'} = \langle \langle \xi^{\mu'} \tanh[\beta(\sum_{\nu=1}^3 \xi^{\nu'} m^{\nu'} - \eta \xi^{4'} m^{4'})] \rangle \rangle_{\xi'}. \quad (43)$$

This is simply the SPEs for the unlearning of ξ^4 . Thus, the stable region of a memory state coincides with one of the stable regions of the memory states for the unlearning of ξ^4 . The same holds true for mixed states. Thus, we next study the case that multiple mixed states are unlearned.

3.2 Case that multiple mixed states are unlearned

In this subsection, we describe numerical results for the unlearning of multiple mixed states. Hereafter, we use new notations to identify solutions. The memory state corresponding

to ξ^μ and the mixed state corresponding to $\xi^{(\mu;\gamma)}$ are denoted by M_μ and $M_{(\mu;\gamma)}$, respectively. In particular, for simplicity, when $\gamma = (1, 1, 1)$, we use ξ^μ and M_μ instead of $\xi^{(\mu;1,1,1)}$ and $M_{(\mu;1,1,1)}$, respectively. Since the S_2 state is considered to be a variant of the Hopfield attractor, we do not distinguish the M_1 and S_2 states and denote both of them by M_1 . We numerically integrated the evolution equations given by Eq. (26) by the RK method starting from memory states or mixed states in order to obtain stable regions of these states. Suppose that we study the overlaps m^1, m^2, \dots, m^v , and let m^ω correspond to the state M_μ or $M_{(\mu;\gamma)}$. For each (T, η) , we set the initial conditions to $m^\omega = 0.9$ and $m^v = \langle \xi^\omega \xi^v \rangle + r$ for $v \neq \omega$, where r is a random number in $[-0.1, 0.1]$. After the trajectory converges, we determine the maximum overlap m^{μ_m} among m^1, \dots, m^v . We judge that the parameter (T, η) is in the stable region of M_μ or $M_{(\mu;\gamma)}$ if $\mu_m = \omega$ and $|m^{\mu_m}| > 0.1$. The criterion of the convergence is that $\|\mathbf{m}(t+1) - \mathbf{m}(t)\| < 10^{-5}$ after a transient, i.e., for $t > 10$. Here, $\|\mathbf{m}\| = \sqrt{\sum_{i=1}^v (m^i)^2}$. This criterion to determine the stable region gives similar results to those in Fig. 5 when ξ^4 is unlearned. In Figs. 8-10, we show the results of unlearning two, three, and all mixed states, respectively. We note that the stable regions of memory states are reduced in size for $\eta > 0$ when multiple mixed states are unlearned. The reason for this is considered to be that the mixed states and memory states are correlated. In order to clarify that this is generally true, we study the case that all mixed states are unlearned with equal weight in the next section for a general value of p .

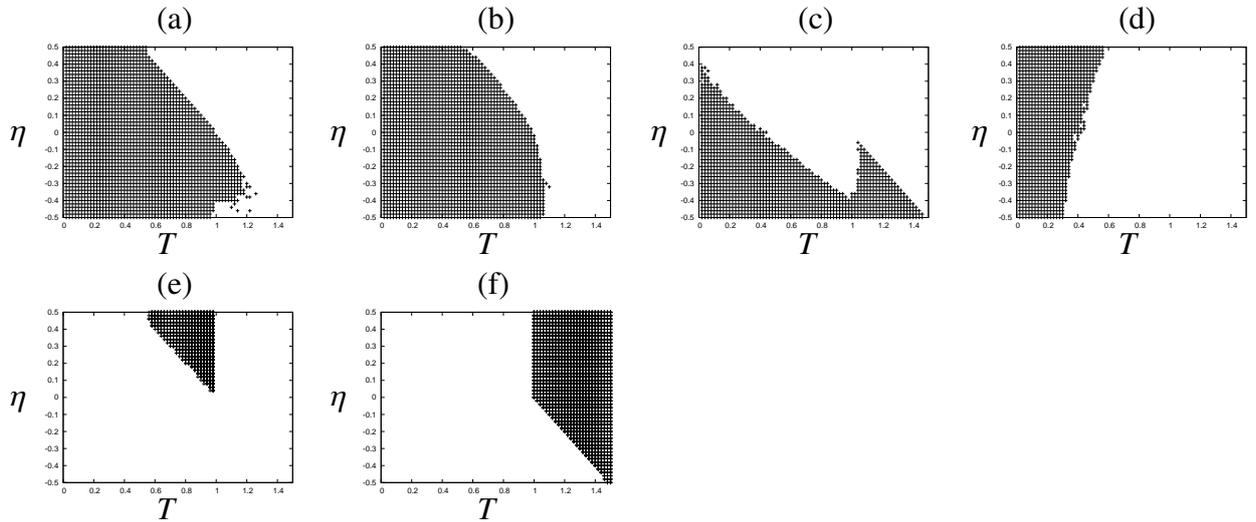


Fig. 8. $p = 3$. $\xi^{(1,2,3)}$ and $\xi^{(1,2,3;1,1,-1)}$ are unlearned with equal weight. Stable regions of solutions are shown in (T, η) space. Dots: stable regions obtained by the RK method. (a) M_1 , (b) M_3 , (c) $M_{(1,2,3)}$, (d) $M_{(1,2,3;1,-1,1)}$, (e) a special state in which $m^1 = -m^2 = m^6, m^3 = m^4 = m^5 = 0$, (f) para state.

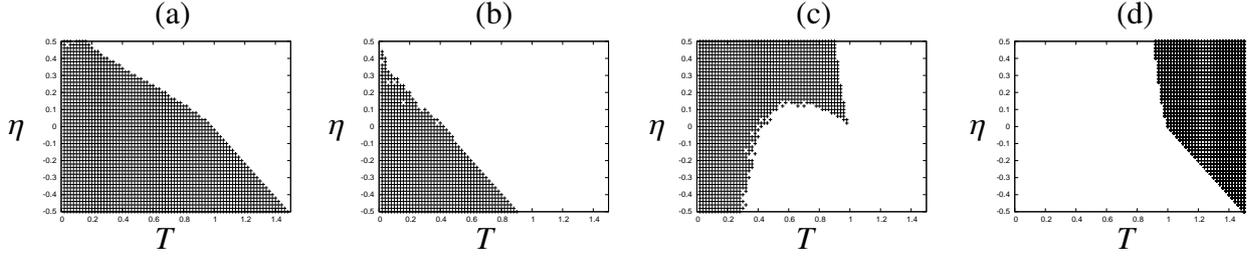


Fig. 9. $p = 3$. Three mixed states $\xi^{(1,2,3)}$, $\xi^{(1,2,3;1,1,-1)}$, and $\xi^{(1,2,3;1,-1,-1)}$ are unlearned with equal weight. Stable regions of solutions are shown in (T, η) space. Dots: stable regions obtained by the RK method. (a) M_1 , (b) $M_{(1,2,3)}$, (c) $M_{(1,2,3;1,-1,-1)}$, (d) para state.

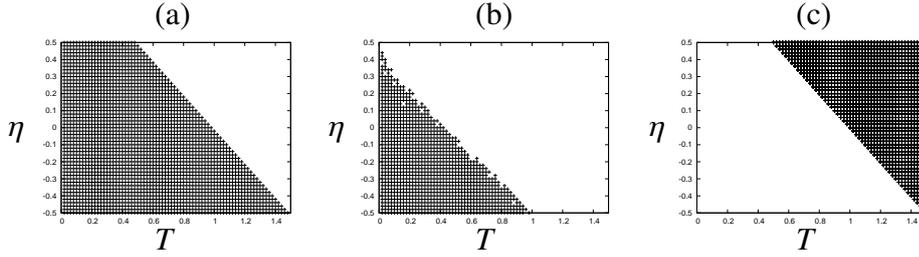


Fig. 10. $p = 3$. All mixed states $\xi^{(1,2,3)}$, $\xi^{(1,2,3;1,1,-1)}$, $\xi^{(1,2,3;1,-1,-1)}$, and $\xi^{(1,2,3;1,-1,-1)}$ are unlearned with equal weight. Stable regions of solutions are shown in (T, η) space. Dots: stable regions obtained by the RK method. (a) M_1 , (b) $M_{(1,2,3)}$, (c) para state.

4. Case that All Mixed States are Unlearned

For convenience, we include $\gamma_1 = -1$ in the summation of mixed states \sum_{ν} in the interaction J_{ij} and rewrite it as

$$J_{ij}^{(U)} = -\frac{\eta}{2N} \sum_{\mathcal{V}'} \xi_i^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \xi_j^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}, \quad (44)$$

where $\mathcal{V}' = \mathcal{V} \cup \{(\boldsymbol{\mu}, \boldsymbol{\gamma}) \text{ with } \gamma_1 = -1\}$. The SPEs are

$$m^{\boldsymbol{\mu}} = \langle\langle \xi^{\boldsymbol{\mu}} \tanh(\beta W) \rangle\rangle, \quad \boldsymbol{\mu} = 1, \dots, p, \quad (45)$$

$$m^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} = \langle\langle \xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} \tanh(\beta W) \rangle\rangle, \quad (\boldsymbol{\mu}; \boldsymbol{\gamma}) \in \mathcal{V}', \quad (46)$$

$$W \equiv \sum_{\nu=1}^p \xi^{\nu} m^{\nu} - \frac{\eta}{2} \sum_{\mathcal{V}'} \xi^{(\boldsymbol{\mu}; \boldsymbol{\gamma})} m^{(\boldsymbol{\mu}; \boldsymbol{\gamma})}. \quad (47)$$

The following relation is useful for deriving the equations below.

$$\text{sgn}(\xi^1 + \xi^2 + \xi^3) = \frac{1}{2}(\xi^1 + \xi^2 + \xi^3 - \xi^1 \xi^2 \xi^3). \quad (48)$$

Let us define the critical coefficients $\eta_c^{\text{pattern}}(T)$ and $\eta_c^{\text{mix}}(T)$ at temperature T below which the pattern and the mixed states exist and are stable as solutions of the SPEs, respectively. These coefficients should satisfy $\eta_c^{\text{pattern}}(T_c^{\text{pattern}}) = 0$ and $\eta_c^{\text{mix}}(T_c^{\text{mix}}) = 0$. Here, $T_c^{\text{pattern}} = 1$ and $T_c^{\text{mix}} \simeq 0.46$ are critical temperatures under which the pattern and mixed states with three

patterns exist for the Hopfield model, respectively. In this section, we derive the formulae for $\eta_c^{\text{pattern}}(0)$ and $\eta_c^{\text{mix}}(0)$. Since all patterns have the same stability and so do all mixed states because of the permutation and inversion symmetries, we take ξ^1 and $\xi^{(\mu_0; \gamma_0)}$ with $\mu_0 = (1, 2, 3)$ and $\gamma_0 = (1, 1, 1)$ as a memory state and a mixed state, respectively. Firstly, let us consider ξ^1 . When $T = 0$, m^1 becomes 1. Thus, we set $s_i = \xi_i^1$. Since $m^\mu = \delta_{\mu,1}$ and $m^{(\mu; \gamma)} = \langle \langle \xi^{(\mu; \gamma)} \xi^1 \rangle \rangle = \frac{1}{2} \gamma_1 \delta_{\mu,1}$, W_i is calculated as

$$W_i = \xi_i^1 - \eta \sum_{1 < \mu_2 < \mu_3} \sum_{\gamma_2} \sum_{\gamma_3} \frac{1}{2} \xi_i^{(\mu; \gamma)} = \left(1 - \frac{\eta}{2}(p-1)(p-2)\right) \xi_i^1. \quad (49)$$

The SPEs given by Eq. (45) become

$$\begin{aligned} \delta_{\mu,1} &= \langle \langle \xi^\mu \tanh(\beta W) \rangle \rangle = \langle \langle \xi^\mu \xi^1 \tanh[\beta \left(1 - \frac{\eta}{2}(p-1)(p-2)\right)] \rangle \rangle \\ &= \delta_{\mu,1} \tanh[\beta \left(1 - \frac{\eta}{2}(p-1)(p-2)\right)]. \end{aligned} \quad (50)$$

Taking the limit $\beta \rightarrow \infty$, if $1 - \frac{\eta}{2}(p-1)(p-2) > 0$, the SPEs are satisfied. Thus, we obtain

$$\eta_c^{\text{pattern}}(0) = \frac{2}{(p-1)(p-2)}. \quad (51)$$

Next, let us study the SPEs given by Eq. (46). They become

$$\frac{1}{2} \gamma_1 \delta_{\mu,1} = \langle \langle \xi^{(\mu; \gamma)} \tanh(\beta W) \rangle \rangle = \langle \langle \xi^{(\mu; \gamma)} \xi^1 \rangle \rangle \tanh[\beta \left(1 - \frac{\eta}{2}(p-1)(p-2)\right)]. \quad (52)$$

This is automatically satisfied in the limit $\beta \rightarrow \infty$ for $\eta < \eta_c^{\text{pattern}}(0)$. Now, let us study the stability of ξ^1 for $\eta < \eta_c^{\text{pattern}}(0)$. In this case, $W \neq 0$. The Hessian matrix is

$$\Lambda_{\mu\nu} = \begin{cases} \delta_{\mu\nu} - \beta \langle \langle \xi^\mu \xi^\nu \cosh^{-2}(\beta W) \rangle \rangle, & \mu, \nu \leq p, \\ \beta \eta \langle \langle \xi^\mu \xi^\nu \cosh^{-2}(\beta W) \rangle \rangle, & \mu = 1, \dots, p, \nu > p, \\ -\eta \delta_{\mu\nu} - \beta \eta^2 \langle \langle \xi^\mu \xi^\nu \cosh^{-2}(\beta W) \rangle \rangle, & \mu, \nu > p. \end{cases} \quad (53)$$

Since $\cosh^{-2}(\beta W) \rightarrow 0$ as $\beta \rightarrow \infty$, the eigenvalues of Λ are 1 (p -fold) and $-\eta$ (u -fold). Thus, ξ^1 is stable for $\eta < \eta_c^{\text{pattern}}(0)$.

Next, let us study the mixed state $\xi^{(\mu_0; \gamma_0)}$. We set $s_i = \xi_i^{(\mu_0; \gamma_0)}$. W_i is calculated as

$$W_i = \frac{1}{2} (\xi_i^1 + \xi_i^2 + \xi_i^3) - \frac{\eta}{2} \sum_{\mu; \gamma} \langle \langle \xi^{(\mu; \gamma)} \xi^{(\mu_0; \gamma_0)} \rangle \rangle \xi_i^{(\mu; \gamma)} \quad (54)$$

$$= \frac{1}{2} (\xi_i^1 + \xi_i^2 + \xi_i^3) - \frac{\eta}{2} (W_i^{(1)} + W_i^{(2)} + W_i^{(3)}), \quad (55)$$

where $W_i^{(1)}$, $W_i^{(2)}$, and $W_i^{(3)}$ are the terms in W_i for which one, two, and all of μ_1, μ_2 , and μ_3 agree with one, two, and all of 1, 2, and 3, respectively. After a little algebra, we obtain

$$W_i^{(1)} = \frac{(p-3)(p-4)}{2} (\xi_i^1 + \xi_i^2 + \xi_i^3), \quad (56)$$

$$W_i^{(2)} = 2(p-3)(\xi_i^1 + \xi_i^2 + \xi_i^3), \quad (57)$$

$$W_i^{(3)} = 2\text{sgn}(\xi_i^1 + \xi_i^2 + \xi_i^3) = 2\xi_i^{\langle \boldsymbol{\mu}_0, \boldsymbol{\gamma}_0 \rangle}. \quad (58)$$

Therefore, we obtain

$$\begin{aligned} W_i &= \frac{1}{2}(\xi_i^1 + \xi_i^2 + \xi_i^3) - \frac{\eta}{4}\left(p(p-3)(\xi_i^1 + \xi_i^2 + \xi_i^3) + 4\xi_i^{\langle \boldsymbol{\mu}_0, \boldsymbol{\gamma}_0 \rangle}\right) \\ &= Vx_i - \eta\xi_i^{\langle \boldsymbol{\mu}_0, \boldsymbol{\gamma}_0 \rangle}, \end{aligned} \quad (59)$$

where $x_i = \xi_i^1 + \xi_i^2 + \xi_i^3$ and $V = \frac{1}{2} - \frac{\eta}{4}p(p-3)$. In the limit of $\beta \rightarrow \infty$, the SPEs given by Eq. (46) become

$$m^{\langle \boldsymbol{\mu}, \boldsymbol{\gamma} \rangle} = \langle \langle \xi^{\langle \boldsymbol{\mu}, \boldsymbol{\gamma} \rangle} \text{sgn}(Vx_i - \eta\xi^{\langle \boldsymbol{\mu}_0, \boldsymbol{\gamma}_0 \rangle}) \rangle \rangle. \quad (60)$$

If $\boldsymbol{\mu} = \boldsymbol{\mu}_0$ and $\boldsymbol{\gamma} = \boldsymbol{\gamma}_0$, we have

$$1 = \langle \langle \text{sgn}(V|x_i| - \eta) \rangle \rangle. \quad (61)$$

From this, we obtain

$$\eta_c^{\text{mix}}(0) = \frac{2}{(p-1)(p-2) + 2}. \quad (62)$$

If none of the elements of $\boldsymbol{\mu}$ agree with the elements of $\boldsymbol{\mu}_0$, both sides of Eq. (60) are 0. Let us consider the case that one of the elements of $\boldsymbol{\mu}$ and $\boldsymbol{\mu}_0$ agree. We assume $\mu_1 = 1$. Then, Eq. (60) becomes

$$\frac{1}{4}\gamma_1 = \langle \langle \text{sgn}(y_i) \text{sgn}(Vx_i - \eta \text{sgn}(x_i)) \rangle \rangle, \quad (63)$$

where $y_i = \gamma_1\xi_i^1 + \gamma_2\xi_i^{\mu_2} + \gamma_3\xi_i^{\mu_3}$. The right-hand side (r.h.s.) of Eq. (63) is calculated as

$$\langle \langle \text{sgn}(y_i) \text{sgn}(Vx_i - \eta \text{sgn}(x_i)) \rangle \rangle = \frac{\gamma_1}{8}(\text{sgn}(3V - \eta) + \text{sgn}(V - \eta)). \quad (64)$$

If $\eta < \eta_c^{\text{mix}}(0)$, $V > \eta$ follows. Thus, the r.h.s. of Eq. (63) is equal to $\frac{\gamma_1}{4}$ and to the left-hand side (l.h.s.) of Eq. (63). Next, we study the case that two elements of $\boldsymbol{\mu}$ and $\boldsymbol{\mu}_0$ agree. We assume $\mu_1 = 1$ and $\mu_2 = 2$. Equation (60) becomes

$$\frac{1}{4}(\gamma_1 + \gamma_2) = \langle \langle \text{sgn}(y_i) \text{sgn}(Vx_i - \eta \text{sgn}(x_i)) \rangle \rangle, \quad (65)$$

where $y_i = \gamma_1\xi_i^1 + \gamma_2\xi_i^2 + \gamma_3\xi_i^{\mu_3}$. We obtain

$$\text{r.h.s. of Eq. (65)} = \frac{1}{8}(\gamma_1 + \gamma_2)(\text{sgn}(3V - \eta) + \text{sgn}(V - \eta)). \quad (66)$$

This is equal to $\frac{1}{4}(\gamma_1 + \gamma_2)$ and to the l.h.s. of Eq. (65) if $\eta < \eta_c^{\text{mix}}(0)$. Finally, let us study the case that $\boldsymbol{\mu} = \boldsymbol{\mu}_0$. Eq. (60) becomes

$$\frac{1}{4}\left(\sum_{k=1}^3 \gamma_k + \prod_{k=1}^3 \gamma_k\right) = \langle \langle \text{sgn}(y_i) \text{sgn}(Vx_i - \eta \text{sgn}(x_i)) \rangle \rangle, \quad (67)$$

where $y_i = \gamma_1 \xi_i^1 + \gamma_2 \xi_i^2 + \gamma_3 \xi_i^3$. Then, we obtain

$$\begin{aligned} \text{r.h.s. of Eq. (67)} = & \frac{1}{8} \left((\gamma_1 + \gamma_2 + \gamma_3 - \gamma_1 \gamma_2 \gamma_3) \text{sgn}(3V - \eta) \right. \\ & \left. + (\gamma_1 + \gamma_2 + \gamma_3 + 3\gamma_1 \gamma_2 \gamma_3) \text{sgn}(V - \eta) \right). \end{aligned} \quad (68)$$

This is equal to $\frac{1}{4}(\gamma_1 + \gamma_2 + \gamma_3 + \gamma_1 \gamma_2 \gamma_3)$ and to the r.h.s. of Eq. (67) if $\eta < \eta_c^{\text{mix}}(0)$. Therefore, all SPEs are satisfied. The stability of $\xi^{(\mu_0; \gamma_0)}$ is the same as that of ξ^1 because the eigenvalues of Λ are 1 (p -fold) and $-\eta$ (u -fold). Thus, all memory states for $\eta < \eta_c^{\text{pattern}}(0)$ and all mixed states for $\eta < \eta_c^{\text{mix}}(0)$ are stable at $T = 0$. From Eqs. (51) and (62), $\eta_c^{\text{pattern}}(0) > \eta_c^{\text{mix}}(0)$ follows. We numerically integrated the evolution equations given by Eq. (26) with $T = 0$ by the RK method starting from the memory states or mixed states for several values of p in order to obtain the stable regions of these states. The criterion of the convergence is that $\|\mathbf{m}(t+1) - \mathbf{m}(t)\| < 0.0001$ after a transient, i.e., for $t > 50$. In Fig. 11, we display the p dependences of $\eta_c^{\text{pattern}}(0)$ and $\eta_c^{\text{mix}}(0)$. The numerical results agree with the theoretical ones reasonably well. As an example, in Fig. 12, for $p = 7$, we display the stable regions for the memory state, mixed state, and para state in the (T, η) plane. Thus, we can eliminate all mixed states

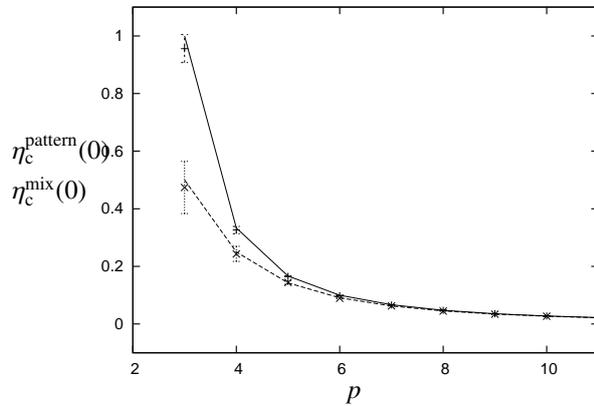


Fig. 11. p dependences of critical unlearning coefficients $\eta_c^{\text{pattern}}(0)$ and $\eta_c^{\text{mix}}(0)$. Curves: theory. Solid curve: $\eta_c^{\text{pattern}}(0)$ and dashed curve: $\eta_c^{\text{mix}}(0)$. Symbols: results obtained by the RK method. Averages are taken from about 50 samples. Vertical lines denote the error bars. +: $\eta_c^{\text{pattern}}(0)$ and \times : $\eta_c^{\text{mix}}(0)$.

and retain all patterns by unlearning in the region \mathcal{R} surrounded by $\eta_c^{\text{mix}}(T)$, $\eta_c^{\text{pattern}}(T)$, $T = 0$, and $\eta = 0$. In order to tune the parameters (T, η) in this region, there are two opposite factors. Since $\eta_c^{\text{pattern}}(0) \rightarrow 0$ and $\eta_c^{\text{mix}}(0) \rightarrow 0$ as $p \rightarrow \infty$, the larger the value of p , the more difficult it is to tune the parameters in \mathcal{R} . The other factor is the better one for tuning. Because $\eta_c^{\text{pattern}}(T_c^{\text{pattern}}) = 0$ and $\eta_c^{\text{mix}}(T_c^{\text{mix}}) = 0$, tuning the parameters in \mathcal{R} seems easier compared with situations in which no such constraints exist.

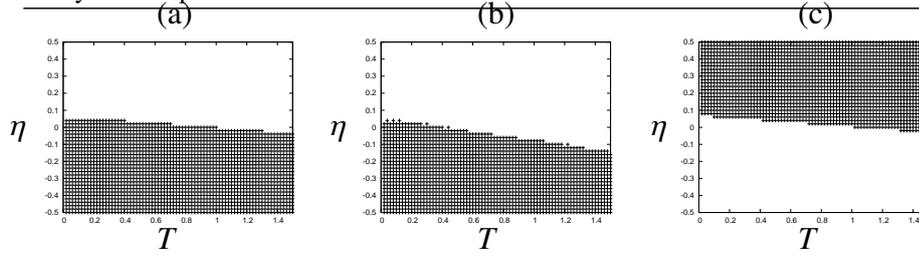


Fig. 12. All mixed states are unlearned. $p = 7$. Stable regions of solutions are shown in (T, η) space obtained by the RK method. (a) M_1 , (b) $M_{(1,2,3)}$, (c) para state.

5. Summary and Discussion

We studied unlearning in the Hopfield model for the case of finite loading of patterns in order to remove mixed states and strengthen the stability of memory states. For the finite loading case, mixed states are the only spurious states. As a method of unlearning, we added terms of mixed states with coefficients ζ_μ to the Hamiltonian of the Hopfield model. We set $\zeta_\mu = -\eta$ throughout the paper when we perform numerical calculations. These terms have the effects of removing the mixed states for $\eta > 0$ and learning them for $\eta < 0$.

In this work, firstly, we studied general situations that any number of mixed states is unlearned and obtained the following results. We derived the SPEs and the evolution equations by introducing sublattices and sublattice overlaps. We postulated that the stability at the equilibrium is determined by the condition that the para state is stable at high temperatures. We proved that the static stability determined by this assumption agrees with the dynamic stability in general situations and that any stable state of the Hopfield model ($\eta = 0$) continuously changes and is statically and dynamically stable for sufficiently small unlearning coefficients. Furthermore, we proved that the stationary solution of the SPEs becomes unstable statically and dynamically at the same time in generic situations.

Next, we studied the case of $p = 3$ in detail. In the case that the single mixed state $\xi^4 = \xi^{(1,1,1)}$ is unlearned, we obtained the following results. For $\eta \neq 0$, the Hopfield attractor does not exist; instead, the S_2 state appears. In the S_2 state, $m^1 (> 0)$ is largest and $m_2 = m_3$. This solution is regarded as a variant of the Hopfield attractor. Similarly, the mixed state $M_5^{(H)}$ for $\eta = 0$, which is characterized by $m^1 = m^2 = -m^3$, changes to M_5 with $m^1 = m^2 \neq -m^3$ for $\eta \neq 0$. On the other hand, the mixed state characterized by $m_1 = m_2 = m_3$ exists for $\eta \neq 0$ as well.

As η increases from 0, the region where the mixed state M_4 is stable decreases in size and disappears at $\eta = 0.5$. On the other hand, the stable S_2 region exists up to $T = 1$ for any $\eta (> 0)$, at least for $\eta \leq 0.5$. That is, by unlearning, the spurious state ξ^4 is deleted while memory states are retained.

The S_3 state in which m_1, m_2 , and m_3 are all different could not be found by numerically solving the SPEs, although in the MCMC simulations, the S_3 solution seemed to exist for small values of N and near the critical temperature T_c with η fixed. It turned out that this phenomenon is a finite size effect because the S_3 state tends to the S_2 state as N becomes large.

In the case of $\eta < 0$, that is, in the case of the learning of the mixed state ξ^4 , as η decreases, the region where the mixed state M_4 is stable increases in size. On the other hand, the region where the S_2 state is stable decreases slowly compared with the decrease in the stable region of M_4 for $\eta > 0$ and disappears at about $\eta \sim -1.73$. We also studied the stability of other mixed states. It is sufficient to study the stability of M_5 owing to the permutation and inversion symmetries that exist in the system. We found that the stable region of M_5 is almost unchanged for $\eta > 0$ and decreases very slowly in size as $|\eta|$ increases for $\eta < 0$. Therefore, we studied the unlearning of other mixed states, $\xi^{(1,2,3;1,1,-1)}$, $\xi^{(1,2,3;1,-1,1)}$, and $\xi^{(1,2,3;1,-1,-1)}$, in order to delete these mixed states. We found that when multiple mixed states are unlearned, the stable regions of memory states are reduced for $\eta > 0$. The reason for this is considered to be that mixed states and memory states are correlated. In order to clarify whether this is generally true, we studied the case that all mixed states are unlearned with equal weight for a general value of p . We defined the critical coefficients $\eta_c^{\text{pattern}}(T)$ and $\eta_c^{\text{mix}}(T)$ at the temperature T below which the pattern and the mixed states exist and are stable as solutions of the SPEs, respectively. We derived the formulae for $\eta_c^{\text{pattern}}(0)$ and $\eta_c^{\text{mix}}(0)$ and their p dependence was numerically confirmed. Therefore, we can eliminate all mixed states and retain all patterns by the unlearning of all mixed states in the region \mathcal{R} surrounded by $\eta_c^{\text{mix}}(T)$, $\eta_c^{\text{pattern}}(T)$, $T = 0$, and $\eta = 0$. In order to tune the parameters (T, η) in this region, there are two opposite factors. Since $\eta_c^{\text{pattern}}(0) \rightarrow 0$ and $\eta_c^{\text{mix}}(0) \rightarrow 0$ as $p \rightarrow \infty$, the larger the value of p , the more difficult it is to tune the parameters in \mathcal{R} . The other factor is the better one for tuning. These coefficients satisfy $\eta_c^{\text{pattern}}(T_c^{\text{pattern}}) = 0$ and $\eta_c^{\text{mix}}(T_c^{\text{mix}}) = 0$, where $T_c^{\text{pattern}} = 1$ and $T_c^{\text{mix}} \simeq 0.46$ are critical temperatures below which the pattern and mixed states with three patterns exist for the Hopfield model, respectively. Owing to these constraints, tuning the parameters in \mathcal{R} seems easier compared with situations in which no such constraints exist.

Appendix A: Evolution Equations for Sublattice Overlaps

In this Appendix, we derive evolution equations for sublattice overlaps $\{\mathcal{M}_l\}$. Let $p_t(\mathcal{M})$ be the probability density of $\mathcal{M} = (\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_2^p)$ at time t ,

$$p_t(\mathcal{M}) = \text{Tr}_{\mathcal{S}} p_t(\mathbf{s}) \prod_{l'=1}^{2^p} \delta(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i).$$

Its time evolution is calculated as

$$\begin{aligned} \frac{\partial}{\partial t} p_t(\mathcal{M}) &= \text{Tr}_{\mathcal{S}} \frac{\partial p_t(\mathbf{s})}{\partial t} \prod_{l'=1}^{2^p} \delta(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i) \\ &= \text{Tr}_{\mathcal{S}} \sum_{l=1}^{2^p} \sum_{k \in \Lambda_l} w_k(\mathbf{s}) p_t(\mathbf{s}) \left[\delta\left\{ \mathcal{M}_l - \frac{2^p}{N} \sum_{i \in \Lambda_l} s_i - \frac{2^p}{N} (s'_k - s_k) \right\} - \delta\left(\mathcal{M}_l - \frac{2^p}{N} \sum_{i \in \Lambda_l} s_i \right) \right] \\ &\quad \times \prod_{l' \neq l} \delta\left(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i \right) \\ &= \text{Tr}_{\mathcal{S}} \sum_{l=1}^{2^p} \sum_{k \in \Lambda_l} w_k(\mathbf{s}) p_t(\mathbf{s}) \left\{ \frac{\partial}{\partial \mathcal{M}_l} \delta\left(\mathcal{M}_l - \frac{2^p}{N} \sum_{i \in \Lambda_l} s_i \right) \cdot \left(-\frac{2^p}{N} \right) (s'_k - s_k) \prod_{l' \neq l} \delta\left(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i \right) \right\} \\ &= \text{Tr}_{\mathcal{S}} \sum_{l=1}^{2^p} \sum_{k \in \Lambda_l} p_t(\mathbf{s}) \frac{2^p}{N} 2w_k(\mathbf{s}) s_k \frac{\partial}{\partial \mathcal{M}_l} \prod_{l'} \delta\left(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i \right), \\ &= \text{Tr}_{\mathcal{S}} \sum_{l=1}^{2^p} \frac{\partial}{\partial \mathcal{M}_l} \left[\left\{ \mathcal{M}_l - \frac{2^p}{N} \sum_{k \in \Lambda_l} \tanh(\beta h_k) \right\} p_t(\mathbf{s}) \prod_{l'} \delta\left(\mathcal{M}_{l'} - \frac{2^p}{N} \sum_{i \in \Lambda_{l'}} s_i \right) \right], \end{aligned}$$

where $s'_k = -s_k$. For $k \in \Lambda_l$, h_k is calculated as

$$\begin{aligned} h_k &= \sum_{j(\neq k)} J_{kj} s_j = \sum_{j(\neq k)} \frac{1}{N} \sum_{\mu=1}^v \zeta_{\mu}^{\xi^{\mu}} \xi_j^{\mu} s_j \simeq \sum_{\mu=1}^v \zeta_{\mu}^{\xi^{\mu}} m^{\mu} = \sum_{\mu=1}^v \zeta_{\mu}^{\xi^{\mu,l}} \frac{1}{2^p} \sum_{l'=1}^{2^p} \xi^{\mu,l'} \mathcal{M}_{l'} \\ &\equiv h^l(\mathcal{M}), \end{aligned}$$

where $\xi^{\mu,l}$ is the value of ξ_k^{μ} for $k \in \Lambda_l$. Therefore, we obtain

$$\sum_{k \in \Lambda_l} \tanh(\beta h_k(\mathcal{M})) = \frac{N}{2^p} \tanh(\beta h^l(\mathcal{M})),$$

$$h^l(\mathcal{M}) = \sum_{\mu=1}^v \zeta_{\mu}^{\xi^{\mu,l}} \frac{1}{2^p} \sum_{l'=1}^{2^p} \xi^{\mu,l'} \mathcal{M}_{l'}.$$

Thus, we obtain

$$\begin{aligned} \frac{\partial}{\partial t} p_t(\mathcal{M}) &= \sum_l \frac{\partial}{\partial \mathcal{M}_l} \left[\left\{ \mathcal{M}_l - \tanh(\beta h^l(\mathcal{M})) \right\} \text{Tr}_{\mathcal{S}} p_t(\mathbf{s}) \prod_l \delta\left(\mathcal{M}_l - \frac{2^p}{N} \sum_{i \in \Lambda_l} s_i \right) \right] \\ &= \sum_l \frac{\partial}{\partial \mathcal{M}_l} \left[\left\{ \mathcal{M}_l - \tanh(\beta h^l(\mathcal{M})) \right\} p_t(\mathcal{M}) \right]. \end{aligned}$$

This implies

$$\frac{d\mathcal{M}_l}{dt} = -\mathcal{M}_l + \tanh[\beta h^l(\mathcal{M})], \quad l = 1, \dots, 2^p. \quad (\text{A}\cdot 1)$$

$h^l(\mathcal{M})$ is rewritten as

$$h^l(\mathcal{M}) = \sum_{\mu=1}^v \zeta_\mu \Xi_{\mu,l} \frac{1}{2^p} (\Xi \mathcal{M})_\mu,$$

where Ξ is the $v \times 2^p$ matrix of which (μ, l) element is $\Xi_{\mu,l} = \xi^{\mu,l}$.

Appendix B: Equivalence of Static and Dynamic Stabilities of Solutions of SPEs

Let us study the stability of the stationary solutions of evolution equations given by Eq. (26). Since the derivatives of f with respect to the overlaps are

$$\frac{\partial f}{\partial m^\mu} = \zeta_\mu m^\mu - \zeta_\mu \langle \xi^\mu \tanh[\beta (\sum_{\nu=1}^v \zeta_\nu m^\nu \xi^\nu)] \rangle, \quad \mu = 1, 2, \dots, v, \quad (\text{B}\cdot 1)$$

Eq. (26) is rewritten as

$$\frac{dm^\mu}{dt} = -\frac{1}{\zeta_\mu} \frac{\partial f}{\partial m^\mu}, \quad \mu = 1, 2, \dots, v. \quad (\text{B}\cdot 2)$$

Let us define \mathbf{m} as

$$\mathbf{m} = (m^1, m^2, \dots, m^v)^T, \quad (\text{B}\cdot 3)$$

where T implies the transpose. Let \mathbf{m}^* be a stationary state of Eq. (26) and $\delta\mathbf{m}$ be a deviation from \mathbf{m}^* ,

$$\mathbf{m}(t) = \mathbf{m}^* + \delta\mathbf{m}, \quad (\text{B}\cdot 4)$$

$$\mathbf{m}^* = (m^{1*}, m^{2*}, \dots, m^{v*})^T, \quad (\text{B}\cdot 5)$$

$$\delta\mathbf{m} = (\delta m^1, \delta m^2, \dots, \delta m^v)^T. \quad (\text{B}\cdot 6)$$

Thus, we obtain

$$\frac{d}{dt} \delta m^\mu = -\frac{1}{\zeta_\mu} \frac{\partial f}{\partial m^\mu} \Big|_{\mathbf{m}^* + \delta\mathbf{m}} \simeq -\frac{1}{\zeta_\mu} \sum_{\nu=1}^v \frac{\partial^2 f}{\partial m^\mu \partial m^\nu} \delta m^\nu, \quad \mu = 1, 2, \dots, v. \quad (\text{B}\cdot 7)$$

By defining $\Lambda_{\mu\nu} = \frac{\partial^2 f}{\partial m^\mu \partial m^\nu} \Big|_{\mathbf{m}=\mathbf{m}^*}$, Eq. (B.7) is rewritten as

$$\frac{d}{dt} \delta\mathbf{m} = -\hat{\Lambda} \delta\mathbf{m}, \quad (\text{B}\cdot 8)$$

where

$$\hat{\Lambda} = \begin{pmatrix} \Lambda_{11} & \cdots & \Lambda_{1p} & \Lambda_{1p+1} & \cdots & \Lambda_{1v} \\ \Lambda_{21} & \cdots & \Lambda_{2p} & \Lambda_{2p+1} & \cdots & \Lambda_{2v} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \Lambda_{p1} & \cdots & \Lambda_{pp} & \Lambda_{pp+1} & \cdots & \Lambda_{pv} \\ \frac{1}{\zeta_{p+1}}\Lambda_{p+11} & \cdots & \frac{1}{\zeta_{p+1}}\Lambda_{p+1p} & \frac{1}{\zeta_{p+1}}\Lambda_{p+1p+1} & \cdots & \frac{1}{\zeta_{p+1}}\Lambda_{p+1v} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{1}{\zeta_v}\Lambda_{v1} & \cdots & \frac{1}{\zeta_v}\Lambda_{vp} & \frac{1}{\zeta_v}\Lambda_{vp+1} & \cdots & \frac{1}{\zeta_v}\Lambda_{vv} \end{pmatrix}. \quad (\text{B}\cdot\text{9})$$

The solution of Eq. (B·8) is

$$\delta\mathbf{m}(t) = e^{-\hat{\Lambda}t}\delta\mathbf{m}(0). \quad (\text{B}\cdot\text{10})$$

Therefore, the stability condition for \mathbf{m}^* is that the real parts of all eigenvalues of $\hat{\Lambda}$ are positive. Let us define a $v \times v$ matrix $G(\{\zeta_\mu\})$ as

$$G(\{\zeta_\mu\}) = G(\{\zeta_\mu\})^T = \begin{pmatrix} \zeta_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \zeta_2 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \ddots & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \zeta_p & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \zeta_{p+1} & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & \ddots & 0 \\ 0 & \cdots & 0 & 0 & 0 & \cdots & \zeta_v \end{pmatrix}. \quad (\text{B}\cdot\text{11})$$

By defining $\delta\tilde{\mathbf{m}} = G(\{\sqrt{|\zeta_\mu|}\})\delta\mathbf{m}$, from Eq. (B·8), we obtain

$$\frac{d}{dt}\delta\tilde{\mathbf{m}} = -\tilde{\Lambda}\delta\tilde{\mathbf{m}}, \quad (\text{B}\cdot\text{12})$$

$$\tilde{\Lambda} = G(\{\frac{\text{sgn}(\zeta_\mu)}{\sqrt{|\zeta_\mu|}}\})^T \Lambda G(\{\frac{1}{\sqrt{|\zeta_\mu|}}\}). \quad (\text{B}\cdot\text{13})$$

If $\zeta_\mu > 0$ for all μ , we obtain

$$\tilde{\Lambda} = G(\{\frac{1}{\sqrt{\zeta_\mu}}\})^T \Lambda G(\{\frac{1}{\sqrt{\zeta_\mu}}\}).$$

Therefore, by Sylvester's law of inertia, the numbers of positive, negative, and zero eigenvalues are the same for $\tilde{\Lambda}$ and Λ . Thus, in this case, the dynamic and static stabilities are the same. On the other hand, if a negative ζ_μ exists, we have to use a different argument. It is necessary to study the Hessian matrix in detail.

Hereafter, the ζ are expressed as follows:

$$\zeta_\mu = 1, \quad \mu = 1, \cdots, p \quad (\text{B}\cdot\text{14})$$

$$\zeta_\mu = -\eta_\mu, \quad \mu = p + 1, \dots, v. \quad (\text{B}\cdot 15)$$

We assume that all the η_μ take nonzero values. We have the following relation:

$$\hat{\Lambda}(\boldsymbol{\eta}) \equiv \hat{\Lambda}(\boldsymbol{\zeta}) = G\left(\left\{-\frac{1}{\eta_\mu}\right\}\right)\Lambda(\boldsymbol{\eta}), \quad (\text{B}\cdot 16)$$

where we define

$$\boldsymbol{\eta} = (\eta_1, \eta_2, \dots, \eta_v), \quad (\text{B}\cdot 17)$$

$$\boldsymbol{\zeta} = (\zeta_1, \zeta_2, \dots, \zeta_v). \quad (\text{B}\cdot 18)$$

Here, we explicitly write down the $\boldsymbol{\eta}$ dependences of Λ and $\hat{\Lambda}$. Let $\lambda_1(\boldsymbol{\eta}), \lambda_2(\boldsymbol{\eta}), \dots, \lambda_v(\boldsymbol{\eta})$ be the eigenvalues of $\Lambda(\boldsymbol{\eta})$. Similarly, let $\hat{\lambda}_1(\boldsymbol{\eta}), \hat{\lambda}_2(\boldsymbol{\eta}), \dots, \hat{\lambda}_v(\boldsymbol{\eta})$ be the eigenvalues of $\hat{\Lambda}(\boldsymbol{\eta})$. Let us study the eigenvalues of $\Lambda(\boldsymbol{\eta})$ for the stationary state \mathbf{m}^* . The components of $\Lambda(\boldsymbol{\eta})$ are expressed as

$$\Lambda(\boldsymbol{\eta})_{\mu\nu} = \begin{cases} \delta_{\mu\nu} - \beta \langle \xi^\mu \xi^\nu \cosh^{-2} W^* \rangle, & \mu, \nu \leq p, \\ \beta \eta_\nu \langle \xi^\mu \xi^\nu \cosh^{-2} W^* \rangle, & \mu = 1, \dots, p, \nu > p, \\ -\eta_\mu \delta_{\mu\nu} - \beta \eta_\mu \eta_\nu \langle \xi^\mu \xi^\nu \cosh^{-2} W^* \rangle, & \mu, \nu > p, \end{cases} \quad (\text{B}\cdot 19)$$

where $W^* \equiv \beta(\sum_{v=1}^3 m^{v*} \xi^v - \sum_{v=p+1}^v \eta_v m^{v*} \xi^v)$. We solve the characteristic equation

$$\left| \Lambda(\boldsymbol{\eta}) - \lambda E_v \right| = 0, \quad (\text{B}\cdot 20)$$

where E_v is the $v \times v$ unit matrix. Since a solution of Eq. (B-20) is a function of $\boldsymbol{\eta}$, we denote it by $\lambda(\boldsymbol{\eta})$. Then, we have

$$h(\boldsymbol{\eta}) \equiv \left| \Lambda(\boldsymbol{\eta}) - \lambda(\boldsymbol{\eta}) E_v \right| = 0. \quad (\text{B}\cdot 21)$$

Assuming that the $|\eta_\mu|$ are small, we expand $\Lambda(\boldsymbol{\eta})$ and $\lambda(\boldsymbol{\eta})$ as Taylor series,

$$\begin{aligned} \Lambda(\boldsymbol{\eta}) &= \Lambda(\mathbf{0}) + \sum_{v(>p)} \frac{\partial \Lambda}{\partial \eta_v}(\mathbf{0}) \eta_v + \dots, \\ \lambda(\boldsymbol{\eta}) &= \lambda(\mathbf{0}) + \sum_{v(>p)} \frac{\partial \lambda}{\partial \eta_v}(\mathbf{0}) \eta_v + \dots. \end{aligned}$$

$\Lambda(\mathbf{0})$ is expressed as

$$\Lambda(\mathbf{0}) = \begin{bmatrix} \Lambda_p^{(H)} & 0_{p,u} \\ 0_{u,p} & 0_{u,u} \end{bmatrix}, \quad (\text{B}\cdot 22)$$

where $0_{l,m}$ is an $l \times m$ zero matrix and $\Lambda_p^{(H)}$ is the Hessian matrix for the Hopfield model,

which is a $p \times p$ matrix. It is expressed as

$$\Lambda_p^{(H)} = \begin{bmatrix} 1 - \beta \langle \langle \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & -\beta \langle \langle \xi^1 \xi^2 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & \cdots & -\beta \langle \langle \xi^1 \xi^p \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle \\ -\beta \langle \langle \xi^2 \xi^1 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & 1 - \beta \langle \langle \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & \cdots & -\beta \langle \langle \xi^2 \xi^p \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle \\ \cdots & \cdots & \cdots & \cdots \\ -\beta \langle \langle \xi^p \xi^1 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & -\beta \langle \langle \xi^p \xi^2 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & \cdots & 1 - \beta \langle \langle \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle \end{bmatrix} \quad (\text{B}\cdot 23)$$

where $W^{(H)} = \beta \sum_{v=1}^p m^{v(H)} \xi^v$ and $m^{v(H)} (v = 1, \dots, p)$ is the stationary state of the Hopfield model. Now, let us consider the situation that all the unlearning coefficients are of the same order and assume $\eta_\mu = \bar{\eta}_\mu \epsilon$ ($\mu = p+1, \dots, v$), where ϵ is a small parameter, $\bar{\eta}_\mu = O(\epsilon^0)$, and $\bar{\eta}_\mu > 0$. The term of $O(\epsilon^0)$ in Eq. (B.21) is $h(\mathbf{0})$, which becomes

$$h(\mathbf{0}) = \left| \Lambda(\mathbf{0}) - \lambda(\mathbf{0}) E_p \right| = \begin{vmatrix} \Lambda_p^{(H)} - \lambda(\mathbf{0}) E_p & 0_{p,u} \\ 0_{u,p} & -\lambda(\mathbf{0}) E_{u,u} \end{vmatrix} \quad (\text{B}\cdot 24)$$

$$= \left| \Lambda_p^{(H)} - \lambda(\mathbf{0}) E_3 \right| (-\lambda(\mathbf{0}))^u = 0. \quad (\text{B}\cdot 25)$$

The solutions of $\left| \Lambda_p^{(H)} - \lambda(\mathbf{0}) E_p \right| = 0$ are the eigenvalues of the stationary state for the Hopfield model, and we denote them by $\lambda_i^{(H)}$ ($i = 1, \dots, p$). Thus, we have $\lambda_i(\mathbf{0}) = \lambda_i^{(H)}$ ($i = 1, \dots, p$). We denote other solutions as $\lambda_\nu(\mathbf{0}) = 0$ ($\nu = p+1, \dots, v$). If the stationary state of the Hopfield model is stable, $\lambda_i^{(H)} > 0$ ($i = 1, \dots, p$). Thus, let us consider $\lambda_\nu(\eta)$ ($\nu > p$). Let us define $h_\nu(\eta)$ as $h(\eta)$ evaluated at $\lambda = \lambda_\nu(\eta)$. Since the components of the $(p+1)$ th to the ν th rows of $\Lambda(\eta)$ are of order ϵ or of order higher than ϵ , the lowest order of $h_\nu(\eta)$ is $O(\epsilon^u)$ and it is given by $\frac{\partial^u h_\nu}{\partial \eta_{p+1} \cdots \partial \eta_\nu}(\mathbf{0}) \bar{\eta}_{p+1} \cdots \bar{\eta}_\nu \epsilon^u$. $\frac{\partial^u h_\nu}{\partial \eta_{p+1} \cdots \partial \eta_\nu}(\mathbf{0})$ is calculated as

$$\frac{\partial^u h_\nu}{\partial \eta_{p+1} \cdots \partial \eta_\nu}(\mathbf{0}) = \quad (\text{B}\cdot 26)$$

$$\begin{vmatrix} \Lambda_p^{(H)} & 0_{p,u} \\ \beta \langle \langle \xi^{p+1} \xi^1 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & \cdots & \beta \langle \langle \xi^{p+1} \xi^p \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & -1 - \frac{\partial \lambda_\nu}{\partial \eta_{p+1}} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \cdots & 0 \\ \beta \langle \langle \xi^\nu \xi^1 \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & \cdots & \beta \langle \langle \xi^\nu \xi^p \frac{1}{\cosh^2 W^{(H)}} \rangle \rangle & 0 & \cdots & 0 & -1 - \frac{\partial \lambda_\nu}{\partial \eta_\nu} \end{vmatrix}$$

$$= \left| \Lambda_p^{(H)} \right| \prod_{\mu=p+1}^v \left(-1 - \frac{\partial \lambda_\nu}{\partial \eta_\mu} \right). \quad (\text{B}\cdot 27)$$

Since $\lambda_i^{(H)} > 0$ ($i = 1, \dots, p$), $\left| \Lambda_p^{(H)} \right| \neq 0$ follows. Thus, we have

$$\frac{\partial \lambda_\nu}{\partial \eta_\mu}(\mathbf{0}) = -1, \quad \mu = p+1, \dots, v. \quad (\text{B}\cdot 28)$$

Therefore, up to $O(\epsilon)$, $\lambda_\nu(\boldsymbol{\eta})$ is expressed as

$$\lambda_\nu(\boldsymbol{\eta}) \simeq \lambda_\nu(\mathbf{0}) + \sum_{\mu=p+1}^v \frac{\partial \lambda_\nu}{\partial \eta_\mu}(\mathbf{0}) \eta_\mu = - \sum_{\mu=p+1}^v \eta_\mu. \quad (\text{B}\cdot\text{29})$$

Next, let us calculate $\hat{\lambda}_i(\mathbf{0})$. We rewrite $\hat{\Lambda}(\boldsymbol{\eta})$ as

$$\hat{\Lambda}(\boldsymbol{\eta}) = \begin{pmatrix} \Lambda_p(\boldsymbol{\eta}) & B(\boldsymbol{\eta}) \\ C(\boldsymbol{\eta}) & D(\boldsymbol{\eta}) \end{pmatrix}, \quad (\text{B}\cdot\text{30})$$

where Λ_p, B, C, D , and G are defined as

$$\left(\Lambda_p(\boldsymbol{\eta}) \right)_{\mu\nu} = \Lambda_{\mu\nu}(\boldsymbol{\eta}), \quad \mu, \nu = 1, \dots, p, \quad (\text{B}\cdot\text{31})$$

$$B(\boldsymbol{\eta}) = \begin{pmatrix} \Lambda_{1p+1}(\boldsymbol{\eta}) & \cdots & \Lambda_{1v}(\boldsymbol{\eta}) \\ \cdots & \cdots & \cdots \\ \Lambda_{pp+1}(\boldsymbol{\eta}) & \cdots & \Lambda_{pv}(\boldsymbol{\eta}) \end{pmatrix}, \quad p \times u \text{ matrix}, \quad (\text{B}\cdot\text{32})$$

$$C(\boldsymbol{\eta}) = G_u\left(\left\{-\frac{1}{\eta_\nu}\right\}\right) \begin{pmatrix} \Lambda_{p+11}(\boldsymbol{\eta}) & \cdots & \Lambda_{p+1p}(\boldsymbol{\eta}) \\ \cdots & \cdots & \cdots \\ \Lambda_{v1}(\boldsymbol{\eta}) & \cdots & \Lambda_{vp}(\boldsymbol{\eta}) \end{pmatrix}, \quad u \times p \text{ matrix}, \quad (\text{B}\cdot\text{33})$$

$$D(\boldsymbol{\eta}) = G_u\left(\left\{-\frac{1}{\eta_\nu}\right\}\right) \begin{pmatrix} \Lambda_{p+1p+1}(\boldsymbol{\eta}) & \cdots & \Lambda_{p+1v}(\boldsymbol{\eta}) \\ \cdots & \cdots & \cdots \\ \Lambda_{vp+1}(\boldsymbol{\eta}) & \cdots & \Lambda_{vv}(\boldsymbol{\eta}) \end{pmatrix}, \quad u \times u \text{ matrix}, \quad (\text{B}\cdot\text{34})$$

$$G_u\left(\left\{-\frac{1}{\eta_\nu}\right\}\right) = \begin{pmatrix} -\frac{1}{\eta_{p+1}} & 0 & \cdots & 0 \\ 0 & \ddots & \cdots & 0 \\ 0 & 0 & \cdots & -\frac{1}{\eta_v} \end{pmatrix}, \quad u \times u \text{ matrix}. \quad (\text{B}\cdot\text{35})$$

Thus, we have

$$\hat{\Lambda}(\mathbf{0}) = \begin{pmatrix} \Lambda_p^{(H)} & 0_{pu} \\ C(\mathbf{0}) & E_u \end{pmatrix}. \quad (\text{B}\cdot\text{36})$$

Now, let us solve the characteristic equation of $\hat{\Lambda}$ for the stationary state studied above,

$$\hat{h}(\boldsymbol{\eta}) \equiv \left| \hat{\Lambda}(\boldsymbol{\eta}) - \hat{\lambda}(\boldsymbol{\eta}) E_v \right| = 0. \quad (\text{B}\cdot\text{37})$$

Let us expand $\hat{\Lambda}(\boldsymbol{\eta})$ and $\hat{\lambda}(\boldsymbol{\eta})$ as Taylor series,

$$\hat{\Lambda}(\boldsymbol{\eta}) = \hat{\Lambda}(\mathbf{0}) + \sum_{\mu>p} \frac{\partial \hat{\Lambda}}{\partial \eta_\mu}(\mathbf{0}) \eta_\mu + \cdots, \quad (\text{B}\cdot\text{38})$$

$$\hat{\lambda}(\boldsymbol{\eta}) = \hat{\lambda}(\mathbf{0}) + \sum_{\mu>p} \frac{\partial \hat{\lambda}}{\partial \eta_\mu}(\mathbf{0}) \eta_\mu + \cdots. \quad (\text{B}\cdot\text{39})$$

The term of $O(\epsilon^0)$ for Eq. (B·37) is

$$\hat{h}(\mathbf{0}) = \begin{vmatrix} \Lambda_p^{(H)} - \hat{\lambda}(\mathbf{0})E_p & 0_{p,u} \\ C(\mathbf{0}) & (1 - \hat{\lambda}(\mathbf{0}))E_u \end{vmatrix} = 0. \quad (\text{B}\cdot 40)$$

From this, we obtain

$$\left| \Lambda_p^{(H)} - \hat{\lambda}(\mathbf{0})E_p \right| (1 - \hat{\lambda}(\mathbf{0}))^u = 0. \quad (\text{B}\cdot 41)$$

Therefore, we have

$$\left| \Lambda_p^{(H)} - \hat{\lambda}(\mathbf{0})E_p \right| = 0, \quad (\text{B}\cdot 42)$$

$$1 - \hat{\lambda}(\mathbf{0}) = 0. \quad (\text{B}\cdot 43)$$

The first equation is the characteristic equation for the Hopfield model, and we have

$$\hat{\lambda}_i(\mathbf{0}) = \lambda_i^H > 0, \quad i = 1, \dots, p. \quad (\text{B}\cdot 44)$$

The second equation gives

$$\hat{\lambda}_\nu(\mathbf{0}) = 1, \quad \nu = p + 1, \dots, v. \quad (\text{B}\cdot 45)$$

Therefore, for sufficiently small $|\epsilon|$, i.e., $|\epsilon| \ll 1$, we have

$$\hat{\lambda}_i(\boldsymbol{\eta}) > 0, \quad i = 1, \dots, v, \quad (\text{B}\cdot 46)$$

and the stationary state is dynamically stable.

Therefore, for sufficiently small ϵ , we have

$$\begin{aligned} \hat{\lambda}_i(\boldsymbol{\eta}) &> 0 \quad i = 1, \dots, v \\ \lambda_i(\boldsymbol{\eta}) &> 0 \quad i = 1, \dots, p \\ \lambda_\nu(\boldsymbol{\eta}) &\simeq - \sum_{\mu=p+1}^v \eta_\mu, \quad i = p + 1, \dots, v. \end{aligned} \quad (\text{B}\cdot 47)$$

Thus, we conclude that irrespective of the sign of η_μ , any stable stationary state of the Hopfield model is statically and dynamically stable for sufficiently small $|\epsilon|$. In particular, this holds for $\eta_\mu = \eta$ for all μ .

Let us consider the breaking of stability when one of the η , say η_i , changes. For any nonzero ϵ , it follows from Eq. (B·16) that

$$\prod_{i=1}^v \hat{\lambda}_i(\boldsymbol{\eta}) = \left[\prod_{j=p+1}^v \left(-\frac{1}{\eta_j}\right) \right] \prod_{i=1}^v \lambda_i(\boldsymbol{\eta}). \quad (\text{B}\cdot 48)$$

Let $\eta_c(T)$ be the value of η_i where the stationary state becomes statically unstable for the first time when $|\eta_i|$ is increased from 0 at a fixed temperature T . Similarly, let $\hat{\eta}_c(T)$ be the value

of η_i where the stationary state becomes dynamically unstable for the first time when $|\eta_i|$ is increased from 0 at a fixed temperature T . Generically, only one of the λ , say λ_j , changes its sign. Therefore, at $\eta_c(T)$, $\lambda_j = 0$. Here, we separately discuss the cases $\eta_i < 0$ and $\eta_i > 0$. For the negative η_i , by Sylvester's law of inertia, only $\hat{\lambda}_j$ becomes 0 at $\eta_i = \eta_c(T)$. Thus, $\eta_c(T) = \hat{\eta}_c(T)$ follows. On the other hand, for the positive η_i , from Eq. (B·48), generically, only one of the $\hat{\lambda}$, say $\hat{\lambda}_k$, changes its sign at $\eta_i = \hat{\eta}_c(T)$, and then $\eta_c(T) = \hat{\eta}_c(T)$ follows. Therefore, the stationary solution becomes unstable statically and dynamically at the same time for positive and negative η_i in generic situations. This property is proved to hold when $\eta_\mu = \eta$ for all μ by a similar argument.

In general, we cannot say anything about whether this unstable solution becomes stable when η_i is further increased because eigenvalues other than λ_j , $\hat{\lambda}_j$, and $\hat{\lambda}_k$ may change their signs.

Appendix C: SPEs in the Case that the Single Mixed State is Unlearned

The SPEs for the solutions are given as follows.

(1) Hopfield attractor H. $m^2 = m^3 = 0$.

$$m^1 = \frac{1}{4} \left(3 \tanh[\beta(m^1 - \eta m^4)] + \tanh[\beta(m^1 + \eta m^4)] \right), \quad (\text{C}\cdot 1)$$

$$m^2 = \frac{1}{4} \left(\tanh[\beta(m^1 - \eta m^4)] - \tanh[\beta(m^1 + \eta m^4)] \right), \quad (\text{C}\cdot 2)$$

$$m^4 = \frac{1}{4} \left(3 \tanh[\beta(m^1 - \eta m^4)] - \tanh[\beta(m^1 + \eta m^4)] \right). \quad (\text{C}\cdot 3)$$

From the condition of $m^2 = 0$, $\eta = 0$ follows. That is, the Hopfield attractor does not exist when the unlearning term exists in the Hamiltonian.

(2) Mixed state M₄, $m_1 = m_2 = m_3$.

$$m^1 = \frac{1}{4} \left(\tanh[\beta(3m^1 - \eta m^4)] + \tanh[\beta(m^1 - \eta m^4)] \right), \quad (\text{C}\cdot 4)$$

$$m^4 = \frac{1}{4} \left(\tanh[\beta(3m^1 - \eta m^4)] + 3 \tanh[\beta(m^1 - \eta m^4)] \right). \quad (\text{C}\cdot 5)$$

(3) Mixed state M₅, $m_1 = m_2 (\simeq -m_3)$.

$$m^1 = \frac{1}{4} \left(\tanh[\beta(2m^1 + m^3 - \eta m^4)] + \tanh[\beta(2m^1 - m^3 - \eta m^4)] \right), \quad (\text{C}\cdot 6)$$

$$m^3 = \frac{1}{4} \left(\tanh[\beta(2m^1 + m^3 - \eta m^4)] - \tanh[\beta(2m^1 - m^3 - \eta m^4)] \right. \\ \left. + 2 \tanh[\beta(m^3 - \eta m^4)] \right), \quad (\text{C}\cdot 7)$$

$$m^4 = \frac{1}{4} \left(\tanh[\beta(2m^1 + m^3 - \eta m^4)] + \tanh[\beta(2m^1 - m^3 - \eta m^4)] + 2\tanh[\beta(m^3 - \eta m^4)] \right). \quad (\text{C}\cdot 8)$$

From these equations, we note that $m^3 \neq -m^1$ unless $\eta = 0$.

(4) S_2 state.

$$m^1 = \frac{1}{4} \left(\tanh[\beta(m^1 + 2m^2 - \eta m^4)] + 2\tanh[\beta(m^1 - \eta m^4)] + \tanh[\beta(m^1 - 2m^2 + \eta m^4)] \right), \quad (\text{C}\cdot 9)$$

$$m^2 = \frac{1}{4} \left(\tanh[\beta(m^1 + 2m^2 - \eta m^4)] - \tanh[\beta(m^1 - 2m^2 + \eta m^4)] \right), \quad (\text{C}\cdot 10)$$

$$m^4 = \frac{1}{4} \left(\tanh[\beta(m^1 + 2m^2 - \eta m^4)] + 2\tanh[\beta(m^1 - \eta m^4)] - \tanh[\beta(m^1 - 2m^2 + \eta m^4)] \right). \quad (\text{C}\cdot 11)$$

Appendix D: Stability Analysis of Solutions of SPEs When One Mixed State is Unlearned

In the general case that f is a function of seven overlaps m^1, \dots, m^7 , we put $\Lambda_{\mu\nu} = \frac{\partial^2 f}{\partial m^\mu \partial m^\nu}$. For $\mu = 1, \dots, 7$ and $\nu = 1, \dots, 7$, we obtain

$$\Lambda_{\mu\nu} = \frac{\partial^2 f}{\partial m^\mu \partial m^\nu} = \zeta_\mu \left(\delta_{\mu\nu} - \beta \zeta_\nu \langle \langle \xi^\mu \xi^\nu \frac{1}{\cosh^2(\beta \sum_{\tau=1}^7 \zeta_\tau m^\tau \xi^\tau)} \rangle \rangle \right). \quad (\text{D}\cdot 1)$$

The characteristic equation of $\Lambda \equiv \{\Lambda_{\mu\nu}\}$ is $|\Lambda - \lambda E| = 0$. The boundary of a stable region for any solution of the SPEs is determined by $|\Lambda| = 0$ because some eigenvalue becomes 0 at the boundary. When only ξ^4 is unlearned, we put $\zeta_1 = \zeta_2 = \zeta_3 = 1$, $\zeta_4 = -\eta$, and $\zeta_5 = \zeta_6 = \zeta_7 = 0$. In this case, the elements of the Hessian matrix are

$$\Lambda_{\mu\nu} = \delta_{\mu\nu} - \beta \langle \langle \xi^\mu \xi^\nu \frac{1}{\cosh^2[\beta(\sum_{\mu=1}^3 m^\mu \xi^\mu - \eta m^4 \xi^4)]} \rangle \rangle, \quad \text{for } \mu, \nu = 1, 2, 3, \quad (\text{D}\cdot 2)$$

$$\Lambda_{\mu 4} = \beta \eta \langle \langle \xi^\mu \xi^4 \frac{1}{\cosh^2[\beta(\sum_{\mu=1}^3 m^\mu \xi^\mu - \eta m^4 \xi^4)]} \rangle \rangle, \quad \text{for } \mu = 1, 2, 3, \quad (\text{D}\cdot 3)$$

$$\Lambda_{44} = -\eta - \beta \eta^2 \langle \langle \frac{1}{\cosh^2[\beta(\sum_{\mu=1}^3 m^\mu \xi^\mu - \eta m^4 \xi^4)]} \rangle \rangle. \quad (\text{D}\cdot 4)$$

We below derive characteristic equations for the solutions of the SPEs.

(1) Hopfield attractor. We have

$$\Lambda_{11} = \Lambda_{22} = \Lambda_{33}, \Lambda_{12} = \Lambda_{13}, \Lambda_{24} = \Lambda_{34}.$$

The characteristic equation is

$$\begin{aligned}
 |\Lambda - \lambda E| &= \{\lambda - (\Lambda_{11} + \Lambda_{12})\}\{\lambda^3 + (\Lambda_{12} - 2\Lambda_{11} - \Lambda_{44})\lambda^2 \\
 &\quad + (\Lambda_{11}^2 - \Lambda_{11}\Lambda_{12} - \Lambda_{44}\Lambda_{12} + 2\Lambda_{11}\Lambda_{44} - 2\Lambda_{12}^2 - 2\Lambda_{24}^2 - \Lambda_{14}^2)\lambda \\
 &\quad + (\Lambda_{12}\Lambda_{11}\Lambda_{44} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{12}\Lambda_{24}\Lambda_{14} + 2\Lambda_{12}^2\Lambda_{44} + 2\Lambda_{24}^2\Lambda_{11} - \Lambda_{14}^2\Lambda_{12} + \Lambda_{14}^2\Lambda_{11})\} = 0.
 \end{aligned} \tag{D·5}$$

(2) P state. We have

$$\Lambda_{11} = \Lambda_{22} = \Lambda_{33} = 1 - \beta, \Lambda_{44} = -\eta - \beta\eta^2, \Lambda_{12} = \Lambda_{13} = \Lambda_{23} = 0, \Lambda_{14} = \Lambda_{24} = \Lambda_{34} = \frac{1}{2}\beta\eta.$$

Thus,

$$|\Lambda - \lambda E| = (\lambda - \Lambda_{11})^2\{\lambda^2 + (-\Lambda_{11} - \Lambda_{44})\lambda + \Lambda_{11}\Lambda_{44} - 3\Lambda_{14}^2\} = 0. \tag{D·6}$$

The solution of Eq. (D·6) is

$$\lambda_1 = \lambda_2 = \Lambda_{11}, \tag{D·7}$$

$$\lambda_{\pm} = \frac{1}{2}\left(\Lambda_{11} + \Lambda_{44} \pm \sqrt{(\Lambda_{11} - \Lambda_{44})^2 + 12\Lambda_{14}^2}\right). \tag{D·8}$$

The stable region for the P state is determined by the conditions that $\lambda_1 \geq \lambda_2 \geq \lambda_3 > 0$ and $\lambda_4 < 0$ for $\eta > 0$, and that all λ are positive for $\eta < 0$. From this, the boundary of the stable P state is given by

$$T = 1 \text{ for } \eta > 0, \tag{D·9}$$

$$\eta = \frac{4T(1 - T)}{4T - 1} \text{ for } \eta < 0. \tag{D·10}$$

(3) Mixed state M₄.

$$\Lambda_{11} = \Lambda_{22} = \Lambda_{33}, \Lambda_{12} = \Lambda_{13} = \Lambda_{23}, \Lambda_{14} = \Lambda_{24} = \Lambda_{34}.$$

$$|\Lambda - \lambda E| = \{\lambda - (\Lambda_{11} - \Lambda_{12})\}^2\{\lambda^2 - (\Lambda_{11} + \Lambda_{44} + 2\Lambda_{12})\lambda + 2\Lambda_{44}\Lambda_{12} + \Lambda_{11}\Lambda_{44} - 3\Lambda_{14}^2\} = 0. \tag{D·11}$$

(4) Mixed state M₅.

$$\Lambda_{11} = \Lambda_{22} = \Lambda_{33}, \Lambda_{13} = \Lambda_{23}, \Lambda_{14} = \Lambda_{24},$$

$$\begin{aligned}
 |\Lambda - \lambda E| &= \{\lambda - (\Lambda_{11} - \Lambda_{12})\}\{\lambda^3 - (2\Lambda_{11} + \Lambda_{44} + \Lambda_{12})\lambda^2 \\
 &\quad + (2\Lambda_{11}\Lambda_{44} + \Lambda_{11}\Lambda_{12} + \Lambda_{11}^2 + \Lambda_{44}\Lambda_{12} - 2\Lambda_{13}^2 - \Lambda_{34}^2 - 2\Lambda_{14}^2)\lambda \\
 &\quad - \Lambda_{11}\Lambda_{44}\Lambda_{12} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{34}\Lambda_{13}\Lambda_{14} + 2\Lambda_{11}\Lambda_{14}^2 + 2\Lambda_{44}\Lambda_{13}^2 + \Lambda_{12}\Lambda_{34}^2 + \Lambda_{11}\Lambda_{34}^2\} = 0.
 \end{aligned}$$

(D·12)

At the boundary, some λ becomes 0. That is,

$$\Lambda_{11} = \Lambda_{12}, \quad (\text{D}\cdot 13)$$

or

$$-\Lambda_{11}\Lambda_{44}\Lambda_{12} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{34}\Lambda_{13}\Lambda_{14} + 2\Lambda_{11}\Lambda_{14}^2 + 2\Lambda_{44}\Lambda_{13}^2 + \Lambda_{12}\Lambda_{34}^2 + \Lambda_{11}\Lambda_{34}^2 = 0. \quad (\text{D}\cdot 14)$$

(5) S_2 state.

$$\Lambda_{11} = \Lambda_{22} = \Lambda_{33}, \Lambda_{12} = \Lambda_{13}, \Lambda_{24} = \Lambda_{34},$$

$$\begin{aligned} |\Lambda - \lambda E| &= \{\lambda - (\Lambda_{11} - \Lambda_{23})\}\{\lambda^3 - (2\Lambda_{11} + \Lambda_{44} + \Lambda_{23})\lambda^2 \\ &\quad + (2\Lambda_{11}\Lambda_{44} + \Lambda_{11}\Lambda_{23} + \Lambda_{44}\Lambda_{23} + \Lambda_{11}^2 - 2\Lambda_{24}^2 - 2\Lambda_{12}^2 - \Lambda_{14}^2)\lambda \\ &\quad - \Lambda_{23}\Lambda_{11}\Lambda_{44} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{12}\Lambda_{24}\Lambda_{14} + 2\Lambda_{11}\Lambda_{24}^2 + 2\Lambda_{44}\Lambda_{12}^2 + \Lambda_{14}^2\Lambda_{23} + \Lambda_{11}\Lambda_{14}^2\} = 0. \end{aligned} \quad (\text{D}\cdot 15)$$

At the boundary, some λ becomes 0. That is,

$$\Lambda_{11} = \Lambda_{23} \quad (\text{D}\cdot 16)$$

or

$$-\Lambda_{23}\Lambda_{11}\Lambda_{44} - \Lambda_{11}^2\Lambda_{44} - 4\Lambda_{12}\Lambda_{24}\Lambda_{14} + 2\Lambda_{11}\Lambda_{24}^2 + 2\Lambda_{44}\Lambda_{12}^2 + \Lambda_{14}^2\Lambda_{23} + \Lambda_{11}\Lambda_{14}^2 = 0. \quad (\text{D}\cdot 17)$$

References

- 1) F. C. Crick, and G. Mitchison, *Nature* **304**, 111 (1983).
- 2) J. J. Hopfield, *Proc. Natl. Acad. Sci. U.S.A.* **79**, 2554 (1982).
- 3) D. J. Amit, H. Gutfreund, and H. Somplinski, *Phys. Rev. A* **32**, 1007 (1985).
- 4) J. Hertz, A. Krogh, and R. G. Palmer, *Introduction to the Theory of Neural Computation* (Addison-Wesley, Redwood City, 1991).
- 5) D. E. Rumelhart, G. E. Hintn, and R. J. Williams, *Learning Internal Representations by Error Propagation, Parallel Distributed Processing, Exploration in the Microstructures of Cognition* (MIT Press, Cambridge, 1986) Vol. 1, Chap. 8, p. 318.
- 6) J. J. Hopfield, D. I. Feinstein, and R. G. Palmer, *Nature* **304**, 158 (1983).
- 7) S. Wimbauer, N. Klemmer, and J. L. van Hemmen, *Neural Networks* **7**, 219 (1994).
- 8) M. C. D. Barrozo, and T. J. P. Penna, *Int. J. Mod. Phys. C* **5**, 503 (1994).
- 9) S. Wimbauer, and J. L. van Hemmen, in *Proceedings on Analysis of Dynamical and Cognitive Systems, Advanced Course* (Springer-Verlag, London, 1995) p. 121.
- 10) J. A. Horas, and E. A. Bea, *Int. J. Neur. Syst.* **12**, 109 (2002).
- 11) K. Nokura, *Phys. Rev. E* **54**, 5571 (1996).
- 12) K. Nokura: *J. Phys. A: Math. Gen.* **31**, 7447 (1998).