

# Dynamics of the Adaptive Natural Gradient Descent Method for Soft Committee Machines

Masato Inoue,<sup>\*</sup> Hyeyoung Park,<sup>†</sup> and Masato Okada<sup>‡</sup>

*Laboratory for Mathematical Neuroscience, RIKEN Brain Science Institute, Saitama 351-0198, Japan*

(Received 13 November 2003; Published 28 May 2004)

Adaptive natural gradient descent (ANGD) realizes natural gradient descent (NGD) without needing to know the input distribution of learning data and reduces the calculation cost from a cubic order to a square order. However, no performance analysis of ANGD has been done. We have developed a statistical-mechanical theory of the simplified version of ANGDPresent address: ft commPresent address: Present address: Present address: Present address: Present address: Present address: ittee machines in on-line learning; this method provides deterministic learning dynamics expressed through a few order parameters, even though ANGD intrinsically holds a large approximated Fisher information matrix. Numerical results obtained using this theory were consistent with those of a simulation, with respect not only to the learning curve but also to the learning failure. Utilizing this method, we numerically evaluated ANGD efficiency and found that ANGD generally performs as well as NGD. We also revealed the key condition affecting the learning plateau in ANGD.

PACS numbers: 02.50.-r, 05.20.-y, 07.05.Mh

## I. INTRODUCTION

Feed-forward multilayer perceptrons are known to have difficulty determining their parameters using a set of training data. This is because of the non-linearity of their activation functions, which prevents the use of analytical estimation methods; e.g., maximum likelihood estimation. An alternative approach is to use the stochastic gradient descent, which introduces an error function for a given learning sample in a supervised learning framework and adjusts the network parameters step by step to reduce the error.

Steepest gradient descent (SGD), equivalent to back-propagation, is a simple and useful gradient descent method, but it suffers from a learning plateau, which is a long learning period with poor error reduction. This learning plateau is caused by the permutation symmetry (i.e., an exchange of two hidden units in the same layer has no effect on the network output) because it creates a saddle structure in the generalization error function [1, 2]. Moreover, the plateau period is prolonged when the weight vectors of hidden units in the teacher network are correlated [3].

In contrast, previous works have shown that natural gradient descent (NGD) [4, 5] has almost optimal learning performance (Fisher efficiency). NGD does not have any plateau if the learning rate is set low enough [6]. NGD is also unaf-

ected by the correlation between the teacher weight vectors [3]. From a general view, one of the greatest advantages of NGD could be its independence of the parameterization of a given network model. Another advantage may be that the premultiplier of the gradient of the error – the inverse of the Fisher information matrix – is not dependent on learning data or the error function and is necessarily positive definite. There are similar methods that use the inverse of the Hessian matrix as the premultiplier. However, these methods may be unstable, because their premultiplier intrinsically depends on learning data and is not necessarily positive definite [7].

Adaptive natural gradient descent (ANGD) [8, 9] is an attractive form of NGD with respect to both the calculation cost and the input distribution of training data. ANGD reduces the calculation cost from a cubic order of the number of network parameters to a square order, and does not need to know the input distribution. Moreover, ANGD retains some of the advantages of NGD; e.g., its approximated inverse of the Fisher information matrix is necessarily positive definite. However, no previous research has quantified the learning performance of ANGD with respect to its practical applicability.

In this paper, we evaluate the learning efficiency of a simplified version of ANGDPresent address: ft commPresent address: Present address: Present address: Present address: Present address: ittee machines in on-line learning. (Soft committee machines consist of simplified two-layer perceptrons.) ANGDPresent address: ft commPresent address: Present address: Present address: Present address: Present address: ittee machines in on-line learning is intrinsically elaborated for on-line learning. On-line learning [7, 10] facilitates analysis because it uses each learning sample only once, so the network state is independent of each learning sample. We employed statistical-mechanical techniques which extract order parameters and make the stochastic learning dynamics converge towards deterministic at the large limit of the input dimension  $N$  [1, 2].

## II. MODEL

We define teacher and student network models, stochastic gradient learning rules of the student parameters (SGD, NGD,

<sup>\*</sup>Present address: Department of Computational Intelligence and Systems Science, Interdisciplinary Graduate School of Science and Engineering, Tokyo Institute of Technology, Yokohama 226-502, Japan; Also at “Intelligent Cooperation and Control”, PRESTO, JST and Department of Otolaryngology, Head and Neck Surgery, Graduate School of Medicine, Kyoto University; Electronic address: minoue@brain.riken.jp

<sup>†</sup>Also at Department of Computer Science, College of Natural Science, Kyungpook National University, Daegu, 702-701, Korea; Electronic address: hypark@knu.ac.kr

<sup>‡</sup>Also at “Intelligent Cooperation and Control”, PRESTO, JST; Electronic address: okada@brain.riken.jp

and ANGD), and an adaptive estimation of the inverse of the Fisher information matrix for ANGD. We use a soft committee machine with  $M$  hidden units as a teacher network and one with  $K$  hidden units as a student network:

$$f_B(\boldsymbol{\xi}) \equiv \sum_{k=1}^M g(\mathbf{B}_k^T \boldsymbol{\xi}), \quad (1)$$

$$f_J(\boldsymbol{\xi}) \equiv \sum_{k=1}^K g(\mathbf{J}_k^T \boldsymbol{\xi}), \quad (2)$$

where  $\bullet^T$  denotes the transpose,  $\mathbf{B}_k \in \mathbb{R}^N$  and  $\mathbf{J}_k \in \mathbb{R}^N$  are column vectors that represent the  $k$ th weight vectors for  $N$ -dimensional input  $\boldsymbol{\xi} \in \mathbb{R}^N$ , and  $g$  denotes the activation function. We define the probability density function for each network for NGD and ANGD. The function for the teacher network is defined with input  $\boldsymbol{\xi}$  and output  $\zeta$  as

$$p_B(\boldsymbol{\xi}, \zeta) \equiv p(\boldsymbol{\xi}) \delta(\zeta - f_B(\boldsymbol{\xi})), \quad (3)$$

where  $p(\boldsymbol{\xi})$  is the input distribution and  $\delta$  is the Dirac delta function. The probability density function for the student network is defined with input  $\boldsymbol{\xi}$  and output  $\zeta'$  using normal distribution  $\mathcal{N}(f_J(\boldsymbol{\xi}), 1)$ :

$$p_J(\boldsymbol{\xi}, \zeta') \equiv p(\boldsymbol{\xi}) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(\zeta' - f_J(\boldsymbol{\xi}))^2\right). \quad (4)$$

The student distribution is modified incrementally by adjusting its parameter vector  $\mathbf{J} \equiv [\mathbf{J}_1^T, \mathbf{J}_2^T, \dots, \mathbf{J}_K^T]^T \in \mathbb{R}^{KN}$  to approximate the teacher distribution.

The error of student output from teacher output is defined as

$$\epsilon_J(\boldsymbol{\xi}) \equiv \frac{1}{2} \{\check{\epsilon}_J(\boldsymbol{\xi})\}^2, \quad (5)$$

where

$$\check{\epsilon}_J(\boldsymbol{\xi}) \equiv f_J(\boldsymbol{\xi}) - f_B(\boldsymbol{\xi}). \quad (6)$$

The generalization error is also defined as the expected error:

$$\epsilon_g(\mathbf{J}) \equiv \langle \epsilon_J(\boldsymbol{\xi}) \rangle_{\boldsymbol{\xi}}, \quad (7)$$

where  $\langle \bullet \rangle_{\circ}$  denotes the expectation of  $\bullet$  with respect to a random variable  $\circ$ .

The parameter update rule in gradient descent can be written, in general, as  $\mathbf{J} := \mathbf{J} + \Delta \mathbf{J}$ , where

$$\Delta \mathbf{J} = -\frac{\eta}{N} \mathbf{M} \nabla_{\mathbf{J}} \epsilon_J(\boldsymbol{\xi}) = -\frac{\eta}{N} \check{\epsilon}_J(\boldsymbol{\xi}) \mathbf{M} \nabla f, \quad (8)$$

while  $\mathbf{M} \in \mathbb{R}^{NK \times NK}$ ,  $\frac{\eta}{N} > 0$  is a scaled update rate,  $\nabla$  denotes the gradient, and  $\nabla f \in \mathbb{R}^{NK}$  denotes  $\nabla_{\mathbf{J}} f_J(\boldsymbol{\xi})$ . We can implement SGD by setting  $\mathbf{M} = \mathbf{I}$  (unit matrix), NGD by setting  $\mathbf{M} = \mathbf{G}^{-1}$  [4], and ANGD by setting  $\mathbf{M} = \hat{\mathbf{G}}^{-1}$  [8, 9]. The  $\mathbf{G}$  denotes the Fisher information matrix of parameter vector  $\mathbf{J}$  defined as

$$\begin{aligned} \mathbf{G} &\equiv \left\langle [\nabla_{\mathbf{J}} \ln p_J(\boldsymbol{\xi}, \zeta')] [\nabla_{\mathbf{J}} \ln p_J(\boldsymbol{\xi}, \zeta')]^T \right\rangle_{\{\boldsymbol{\xi}, \zeta'\}} \\ &= \left\langle [\nabla f] [\nabla f]^T \right\rangle_{\boldsymbol{\xi}}, \end{aligned} \quad (9)$$

while  $\hat{\mathbf{G}}$  is an adaptively approximated matrix of  $\mathbf{G}$  obtained by ANGD. The update rule of  $\hat{\mathbf{G}}$  is given as  $\hat{\mathbf{G}} := \hat{\mathbf{G}} + \Delta \hat{\mathbf{G}}$ , where

$$\Delta \hat{\mathbf{G}} \equiv \frac{\rho}{N} \left[ -\hat{\mathbf{G}} + [\nabla f] [\nabla f]^T \right], \quad (10)$$

while  $0 < \frac{\rho}{N} < 1$  is an update rate. ANGD does not use the input distribution  $p(\boldsymbol{\xi})$ , but shares each input sample  $\boldsymbol{\xi}$  with the update rule of  $\mathbf{J}$ , and approximates  $\mathbf{G}$  step by step. Realistically, rather than Eq. (10), ANGD adopts an exactly equivalent rule using the Sherman-Morrison formula:  $\hat{\mathbf{G}}^{-1} := \hat{\mathbf{G}}^{-1} + \Delta \hat{\mathbf{G}}^{-1}$ , where

$$\Delta \hat{\mathbf{G}}^{-1} = \frac{\frac{\rho}{N}}{1 - \frac{\rho}{N}} \left[ \hat{\mathbf{G}}^{-1} - \frac{\hat{\mathbf{G}}^{-1} [\nabla f] [\nabla f]^T \hat{\mathbf{G}}^{-1}}{1 + \frac{\rho}{N} \{[\nabla f]^T \hat{\mathbf{G}}^{-1} [\nabla f] - 1\}} \right]. \quad (11)$$

Equation (11) offers the great advantage that we can omit the expensive matrix inversion ( $O(\{NK\}^3)$ ) and achieve lower calculation cost  $O(\{NK\}^2)$ . Here,  $\hat{\mathbf{G}}^{-1}$  is always positive definite and symmetric if the initial value is positive definite and symmetric. For an initial value, we choose  $\mathbf{I}$  for simplicity. ANGD is always applicable if SGD is applicable, because  $\nabla f$  is required even in the case of SGD, whereas NGD cannot always be applied when the input distribution is unknown.

When  $\rho \ll 1$ , Eq. (11) can be reduced to a simple form:

$$\Delta \hat{\mathbf{G}}^{-1} = \frac{\rho}{N} \left[ \hat{\mathbf{G}}^{-1} - \hat{\mathbf{G}}^{-1} [\nabla f] [\nabla f]^T \hat{\mathbf{G}}^{-1} \right]. \quad (12)$$

This approximation was introduced by Amari et al [8]. In this paper, we investigate this simplified version of ANGD in detail under the assumption of small  $\rho$ . We also elucidate what happens when this assumption is violated.

### III. THEORETICAL RESULTS

In this section, we show the order parameter expression of the system dynamics in ANGD, where we use both the usual and newly introduced order parameters. With respect to SGD and NGD, the usual order parameters are sufficient to explain the system state because the system has rotation invariance under the assumption of Gaussian input ( $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ); i.e., the system is equivalent to one with rotated weight vectors  $\mathbf{J}_i, \mathbf{B}_i$ . The usual order parameters can also describe the Fisher information matrix  $\mathbf{G}$  and  $\mathbf{G}^{-1}$ . We need new order parameters for ANGD, though, to describe the approximated inverse of the Fisher information matrix  $\hat{\mathbf{G}}^{-1}$ . To make the present paper self-contained, we first briefly summarize the derivation of the usual order parameter equations of the soft committee machines for SGD [1, 2] and NGD [5, 6]. Then, we explain the new order parameter expression for ANGD.

#### A. Generalization error

The usual order parameters are the inner products among all weight vectors:

$$\begin{aligned} \mathbf{Q} &\equiv [Q_{ij}]_{i,j=1,\dots,K}, & Q_{ij} &\equiv \mathbf{J}_i^T \mathbf{J}_j, \\ \mathbf{R} &\equiv [R_{ij}]_{i=1,\dots,K, j=1,\dots,M}, & R_{ij} &\equiv \mathbf{J}_i^T \mathbf{B}_j, \\ \mathbf{T} &\equiv [T_{ij}]_{i,j=1,\dots,M}, & T_{ij} &\equiv \mathbf{B}_i^T \mathbf{B}_j. \end{aligned} \quad (13)$$

Here,  $\mathbf{Q} \in \mathbb{R}^{K \times K}$  means the inner products matrix for the student weight vectors, while  $\mathbf{R} \in \mathbb{R}^{K \times M}$  means the matrix containing the inner products between the student and teacher weight vectors. The  $\mathbf{Q}$  and  $\mathbf{R}$  are updated according to the updating of  $\mathbf{J}$ . Here,  $\mathbf{T} \in \mathbb{R}^{M \times M}$  means the inner products among the teacher weight vectors; this is fixed. The square length of each input, and the inner products between the input and the weight vectors are temporarily used to describe the micro dynamics:

$$\begin{aligned} \chi &\equiv \xi^T \xi, \\ \mathbf{x} &\equiv [x_i]_{i=1, \dots, K}, x_i \equiv \mathbf{J}_i^T \xi, \\ \mathbf{y} &\equiv [y_i]_{i=1, \dots, M}, y_i \equiv \mathbf{B}_i^T \xi. \end{aligned} \quad (14)$$

Here,  $\chi \in \mathbb{R}$  stochastically converges to  $N$  ( $\chi \xrightarrow{P} N$ ) at the large limit of  $N$ . Also,  $\mathbf{x} \in \mathbb{R}^K$  and  $\mathbf{y} \in \mathbb{R}^M$  are random vectors dependent on input  $\xi$ . The distribution of  $\mathbf{z} \equiv [\mathbf{x}^T, \mathbf{y}^T]^T \in \mathbb{R}^{K+M}$  is determined using the order parameters:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$  where

$$\mathbf{C} \equiv \begin{bmatrix} \mathbf{Q} & \mathbf{R} \\ \mathbf{R}^T & \mathbf{T} \end{bmatrix} \in \mathbb{R}^{(K+M) \times (K+M)}. \quad (15)$$

We can then substitute the order parameters and these random variables for all  $N$ -dimensional vectors  $\mathbf{B}_i, \mathbf{J}_i$ , and  $\xi$ . The number of order parameters is sufficiently small because it does not depend on  $N$ , and this facilitates system analysis. For example, the generalization error can be expressed as

$$\epsilon_g(\mathbf{C}) = \frac{1}{2} \int d\mathbf{z} p(\mathbf{z}) \left\{ \sum_{k=1}^{K+M} c_k g(z_k) \right\}^2, \quad (16)$$

where  $c_k$  is 1 if  $k \leq K$  or  $-1$  if  $k > K$ . Here and hereafter, we assume  $g(x) = \text{erf}(x/\sqrt{2})$ , where  $\text{erf}(x) \equiv \frac{2}{\sqrt{\pi}} \int_0^x dt e^{-t^2}$  is the conventional error function. The generalization error can then be rewritten as an analytical function [1],

$$\epsilon_g(\mathbf{C}) = \frac{1}{\pi} \sum_{i,j=1}^{K+M} c_i c_j \arcsin \frac{C_{ij}}{\sqrt{(C_{ii}+1)(C_{jj}+1)}}. \quad (17)$$

### B. Steepest Gradient Descent

The dynamics of the order parameters for SGD can be expressed using the order parameters themselves [1, 2]. From Eq. (8), the update rule for parameter  $\mathbf{J}$  is

$$\Delta \mathbf{J}_i = -\frac{\eta}{N} \delta_i \xi, \quad (18)$$

where

$$\delta \equiv [\delta_i]_{i=1, \dots, K}, \delta_i \equiv \xi_J(\mathbf{z}) g'(x_i), \quad (19)$$

while  $g'$  is the derivative of  $g$ . Thus, the update rules of the order parameters can be written as

$$\begin{aligned} \Delta R_{ij} &= [\Delta \mathbf{J}_i]^T \mathbf{B}_j = -\frac{\eta}{N} \delta_i y_j, \\ \Delta Q_{ij} &= [\mathbf{J}_i + \Delta \mathbf{J}_i]^T [\mathbf{J}_j + \Delta \mathbf{J}_j] - \mathbf{J}_i^T \mathbf{J}_j \\ &= \Delta S_{ij} + \Delta S_{ji} + \Delta \phi_{ij}, \end{aligned} \quad (20)$$

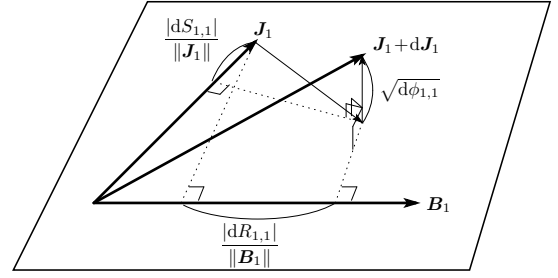


FIG. 1: Intuitive schema of the learning dynamics at the large limit of  $N$ . In the case of  $K = M = 1$ ,  $\mathbf{J}_1$  mainly moves in the current subspace made by  $\mathbf{J}_1$  and  $\mathbf{B}_1$  according to the gradient of the error, but necessarily moves out of this subspace into the null space, or the complementary subspace, of  $\mathbf{J}_1$  and  $\mathbf{B}_1$ , because the small fluctuation of  $\mathbf{J}_1$  in each dimension is summed up resulting in the non-zero term  $\phi_{1,1}$  in Eq. (25).

where

$$\begin{aligned} \Delta S_{ij} &\equiv [\Delta \mathbf{J}_i]^T \mathbf{J}_j = -\frac{\eta}{N} \delta_i x_j, \\ \Delta \phi_{ij} &\equiv [\Delta \mathbf{J}_i]^T \Delta \mathbf{J}_j = \frac{\eta^2 \chi}{N^2} \delta_i \delta_j. \end{aligned} \quad (21)$$

Here and hereafter, we use  $\mathbf{S} \equiv [S_{ij}]_{i,j=1, \dots, K} \in \mathbb{R}^{K \times K}$  and  $\boldsymbol{\phi} \equiv [\phi_{ij}]_{i,j=1, \dots, K} \in \mathbb{R}^{K \times K}$  for convenience.

Next, we introduce time  $\alpha$ , and specify that time  $1/N$  is consumed by each update. At the large limit of  $N$ , the dynamics of the order parameters become continuous and deterministic. For example, the value of  $R_{ij}$  over an infinitesimal time interval ( $d\alpha$ ) after current time  $\alpha$  is not a random variable, although each  $\Delta R_{ij}$  is a random variable.

$$\begin{aligned} R_{ij}^{(\alpha+d\alpha)} &= R_{ij}^{(\alpha)} + \lim_{N \rightarrow \infty} \sum_{\mu=0}^{Nd\alpha-1} \Delta R_{ij}^{(\alpha+\mu/N)} \\ &= R_{ij}^{(\alpha)} - \eta \langle \delta_i y_j \rangle_z d\alpha. \end{aligned} \quad (22)$$

Therefore, the time derivation of  $R_{ij}$  is

$$\frac{dR_{ij}^{(\alpha)}}{d\alpha} = -\eta \langle \delta_i y_j \rangle_z. \quad (23)$$

This expectation with respect to  $\mathbf{z}$  can be solved analytically [1]. The dynamics of  $Q_{ij}$  can be determined similarly, and we get

$$\frac{dQ_{ij}}{d\alpha} = \frac{dS_{ij}}{d\alpha} + \frac{dS_{ji}}{d\alpha} + \frac{d\phi_{ij}}{d\alpha}, \quad (24)$$

where

$$\begin{aligned} \frac{dS_{ij}}{d\alpha} &= -\eta \langle \delta_i x_j \rangle_z, \\ \frac{d\phi_{ij}}{d\alpha} &= \eta^2 \langle \delta_i \delta_j \rangle_z. \end{aligned} \quad (25)$$

The order parameter dynamics suggest that, at the large limit of  $N$ ,  $\mathbf{J}_i$  necessarily moves out of the current direct sum

subspace made by all the weight vectors  $\mathbf{J}_1, \dots, \mathbf{J}_K, \mathbf{B}_1, \dots, \mathbf{B}_M$ , although it mainly moves toward the subspace made by  $\mathbf{B}_1, \dots, \mathbf{B}_M$  (Fig. 1). The direction of this orthogonal movement to that subspace is always chosen randomly, while the time derivative of the square distance of this movement is represented as  $\frac{d\phi_{ij}}{d\alpha}$  in Eq. (24). This orthogonal movement also appears in NGD and ANGd as the  $\frac{d\phi_{ij}}{d\alpha}$  term included in the dynamics of  $\mathcal{Q}_{ij}$ , and, moreover, in the  $\hat{\mathbf{G}}^{-1}$  dynamics for ANGd.

### C. Natural Gradient Descent

The dynamics for NGD can also be expressed through the order parameters. In this subsection, we first determine the Fisher information matrix,  $\mathbf{G}$ , and its inverse,  $\mathbf{G}^{-1}$ , by using only  $\mathcal{Q}$ . Then we derive the motion equations of the order parameters.

Generally, for any  $\mathbf{G}_{ij} \in \mathbb{R}^{K \times K}$ , the  $(i, j)$  block of  $\mathbf{G}$ , can be expressed with both the student weight vectors,  $\mathbf{J}' \equiv [\mathbf{J}_1, \dots, \mathbf{J}_K] \in \mathbb{R}^{N \times K}$ , and all the orthonormal bases of the null

space of the student weight vectors,  $\mathbf{V} \in \mathbb{R}^{N \times (N-K)}$ , which satisfies  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$  and  $\mathbf{J}'^T \mathbf{V} = \mathbf{0}$ , as

$$\mathbf{G}_{ij} = \lambda_{ij} \mathbf{I} + [\mathbf{J}', \mathbf{V}] \begin{bmatrix} \mathbf{\Lambda}_{ij} & \mathbf{\Lambda}'_{ij} \\ \mathbf{\Lambda}''_{ij} & \mathbf{\Lambda}'''_{ij} \end{bmatrix} [\mathbf{J}', \mathbf{V}]^T, \quad (26)$$

where  $\lambda_{ij} \in \mathbb{R}$ ,  $\mathbf{\Lambda}_{ij} \in \mathbb{R}^{K \times K}$ ,  $\mathbf{\Lambda}'_{ij} \in \mathbb{R}^{K \times (N-K)}$ ,  $\mathbf{\Lambda}''_{ij} \in \mathbb{R}^{(N-K) \times K}$ , and  $\mathbf{\Lambda}'''_{ij} \in \mathbb{R}^{(N-K) \times (N-K)}$  (to be exact,  $\lambda_{ij}$  is verbose). We can then rewrite Eq. (26) as,

$$\mathbf{G}_{ij} = \lambda_{ij} \mathbf{I} + \mathbf{J}' \mathbf{\Lambda}_{ij} \mathbf{J}'^T + \mathbf{E}_{ij}, \quad (27)$$

$$\mathbf{E}_{ij} \equiv [\mathbf{J}', \mathbf{V}] \begin{bmatrix} \mathbf{0} & \mathbf{\Lambda}'_{ij} \\ \mathbf{\Lambda}''_{ij} & \mathbf{\Lambda}'''_{ij} \end{bmatrix} [\mathbf{J}', \mathbf{V}]^T, \quad (28)$$

where  $\mathbf{E}_{ij} \in \mathbb{R}^{N \times N}$ .

We next prove  $\mathbf{E}_{ij} = \mathbf{0}$ ; in other words,  $\lambda_{ij}$  and  $\mathbf{\Lambda}_{ij}$  are sufficient parameters to express  $\mathbf{G}$ . By multiplying the identity matrix,  $[\mathbf{J}', \mathbf{V}] [\mathbf{J}', \mathbf{V}]^T [\mathbf{J}', \mathbf{V}]^{-1} [\mathbf{J}', \mathbf{V}]^T$ , we can rewrite  $\mathbf{G}_{ij}$  as

$$\begin{aligned} \mathbf{G}_{ij} &= \langle [\nabla f_i] [\nabla f_j]^T \rangle_{\xi} \\ &= [\mathbf{J}', \mathbf{V}] [\mathbf{J}', \mathbf{V}]^T [\mathbf{J}', \mathbf{V}]^{-1} [\mathbf{J}', \mathbf{V}]^T \langle [\nabla f_i] [\nabla f_j]^T \rangle_{\xi} [\mathbf{J}', \mathbf{V}] [\mathbf{J}', \mathbf{V}]^T [\mathbf{J}', \mathbf{V}]^{-1} [\mathbf{J}', \mathbf{V}]^T \\ &= [\mathbf{J}', \mathbf{V}] \begin{bmatrix} \mathbf{Q}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \left\langle g'(x_i) g'(x_j) \begin{bmatrix} \mathbf{x} \mathbf{x}^T & \mathbf{x} \mathbf{v}^T \\ \mathbf{v} \mathbf{x}^T & \mathbf{v} \mathbf{v}^T \end{bmatrix} \right\rangle_{\xi} \begin{bmatrix} \mathbf{Q}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} [\mathbf{J}', \mathbf{V}]^T \\ &= [\mathbf{J}', \mathbf{V}] \begin{bmatrix} \mathbf{Q}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \left\langle g'(x_i) g'(x_j) \begin{bmatrix} \mathbf{x} \mathbf{x}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \right\rangle_x \begin{bmatrix} \mathbf{Q}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} [\mathbf{J}', \mathbf{V}]^T \\ &= [\mathbf{J}', \mathbf{V}] \begin{bmatrix} \mathbf{Q}^{-1} \langle g'(x_i) g'(x_j) \mathbf{x} \mathbf{x}^T \rangle_x \mathbf{Q}^{-1} & \mathbf{0} \\ \mathbf{0} & \langle g'(x_i) g'(x_j) \rangle_x \mathbf{I} \end{bmatrix} [\mathbf{J}', \mathbf{V}]^T \\ &= \langle g'(x_i) g'(x_j) \rangle_x \mathbf{V} \mathbf{V}^T + \mathbf{J}' \begin{bmatrix} \mathbf{Q}^{-1} \langle g'(x_i) g'(x_j) \mathbf{x} \mathbf{x}^T \rangle_x \mathbf{Q}^{-1} \\ \mathbf{0} \end{bmatrix} \mathbf{J}'^T \\ &= \langle g'(x_i) g'(x_j) \rangle_x \mathbf{I} + \mathbf{J}' \begin{bmatrix} \mathbf{Q}^{-1} \langle g'(x_i) g'(x_j) \mathbf{x} \mathbf{x}^T \rangle_x \mathbf{Q}^{-1} - \langle g'(x_i) g'(x_j) \rangle_x \mathbf{Q}^{-1} \\ \mathbf{0} \end{bmatrix} \mathbf{J}'^T, \end{aligned} \quad (29)$$

where  $\nabla f_i \equiv g'(x_i) \xi$ , while  $\mathbf{v} \equiv \mathbf{V}^T \xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . We used  $\langle g'(x_i) g'(x_j) \mathbf{x} \mathbf{v}^T \rangle_{\xi} = \langle g'(x_i) g'(x_j) \mathbf{x} \rangle_x \langle \mathbf{v}^T \rangle_{\mathbf{v}} = \mathbf{0}$ , and  $\langle \mathbf{v} \mathbf{v}^T \rangle_{\xi} = \mathbf{I}$ . Thus, we can prove  $\mathbf{E}_{ij} = \mathbf{0}$  by letting

$$\lambda_{ij} \equiv \langle g'(x_i) g'(x_j) \rangle_x = \frac{2}{\pi} \left| \begin{array}{cc} \mathcal{Q}_{ii+1} & \mathcal{Q}_{ij} \\ \mathcal{Q}_{ij} & \mathcal{Q}_{jj+1} \end{array} \right|^{-\frac{1}{2}}, \quad (30)$$

$$\begin{aligned} \mathbf{\Lambda}_{ij} &\equiv \mathbf{Q}^{-1} \langle g'(x_i) g'(x_j) \mathbf{x} \mathbf{x}^T \rangle_x \mathbf{Q}^{-1} - \lambda_{ij} \mathbf{Q}^{-1} \\ &= -\lambda_{ij} [\mathbf{e}_i, \mathbf{e}_j] \begin{bmatrix} \mathcal{Q}_{ii+1} & \mathcal{Q}_{ij} \\ \mathcal{Q}_{ij} & \mathcal{Q}_{jj+1} \end{bmatrix}^{-1} [\mathbf{e}_i, \mathbf{e}_j]^T, \end{aligned} \quad (31)$$

where  $|\bullet|$  denotes a determinant, while  $\mathbf{e}_i \in \mathbb{R}^K$  is a unit vector whose  $i$ th element is 1. The  $\lambda_{ij}$  and  $\mathbf{\Lambda}_{ij}$  are expressed with certainty by  $\mathcal{Q}$ .

Next, we determine  $\mathbf{G}^{-1}$  using a similar style to Eq. (27):

$$[\mathbf{G}^{-1}]_{ij} = \theta_{ij} \mathbf{I} + \mathbf{J}' \mathbf{\Theta}_{ij} \mathbf{J}'^T, \quad (32)$$

where  $\theta_{ij} \in \mathbb{R}$  and  $\mathbf{\Theta}_{ij} \in \mathbb{R}^{K \times K}$  correspond to  $\lambda_{ij}$  and  $\mathbf{\Lambda}_{ij}$ , respectively. If we temporarily adopt normalized expressions of  $\mathbf{\Lambda}_{ij}$  and  $\mathbf{\Theta}_{ij}$ , defined as

$$\tilde{\mathbf{\Lambda}}_{ij} \equiv \mathbf{Q}^{\frac{1}{2}} \mathbf{\Lambda}_{ij} \mathbf{Q}^{\frac{1}{2}} + \lambda_{ij} \mathbf{I}, \quad (33)$$

$$\tilde{\mathbf{\Theta}}_{ij} \equiv \mathbf{Q}^{\frac{1}{2}} \mathbf{\Theta}_{ij} \mathbf{Q}^{\frac{1}{2}} + \theta_{ij} \mathbf{I}, \quad (34)$$

matrix multiplication will be simplified; e.g.,

$$\mathbf{G}_{ij} \mathbf{G}_{kl} = [\lambda_{ij} [\mathbf{I} - \mathbf{U} \mathbf{Q}^{-1} \mathbf{U}^T] + \mathbf{U} \mathbf{Q}^{-\frac{1}{2}} \tilde{\mathbf{\Lambda}}_{ij} \mathbf{Q}^{-\frac{1}{2}} \mathbf{U}^T]$$

$$\begin{aligned}
& \times [\lambda_{kl}[\mathbf{I} - \mathbf{U}\mathbf{Q}^{-1}\mathbf{U}^T] + \mathbf{U}\mathbf{Q}^{-\frac{1}{2}}\tilde{\Lambda}_{kl}\mathbf{Q}^{-\frac{1}{2}}\mathbf{U}^T] \\
& = \lambda_{ij}\lambda_{kl}[\mathbf{I} - \mathbf{U}\mathbf{Q}^{-1}\mathbf{U}^T] + \mathbf{U}\mathbf{Q}^{-\frac{1}{2}}\tilde{\Lambda}_{ij}\tilde{\Lambda}_{kl}\mathbf{Q}^{-\frac{1}{2}}\mathbf{U}^T.
\end{aligned} \tag{35}$$

Utilizing this normalization, we can easily obtain

$$\boldsymbol{\theta} = \boldsymbol{\lambda}^{-1}, \tilde{\boldsymbol{\Theta}} = \tilde{\Lambda}^{-1}, \tag{36}$$

where  $\boldsymbol{\lambda} \equiv [\lambda_{ij}]_{i,j=1,\dots,K}$  and  $\boldsymbol{\theta} \equiv [\theta_{ij}]_{i,j=1,\dots,K}$  are  $\mathbb{R}^{K \times K}$  symmetric matrices, while  $\Lambda \equiv [\Lambda_{ij}]_{i,j=1,\dots,K}$ ,  $\tilde{\Lambda} \equiv [\tilde{\Lambda}_{ij}]_{i,j=1,\dots,K}$ ,  $\Theta \equiv [\Theta_{ij}]_{i,j=1,\dots,K}$ , and  $\tilde{\Theta} \equiv [\tilde{\Theta}_{ij}]_{i,j=1,\dots,K}$  are  $\mathbb{R}^{K^2 \times K^2}$  symmetric matrices.

Finally, we obtain the order parameter dynamics for NGD. In NGD, the rule for updating  $\mathbf{J}$  is given by

$$\Delta \mathbf{J}_i = -\frac{\eta}{N} \sum_{k=1}^K \delta_k \mathbf{G}_{ik}^{-1} \boldsymbol{\xi}. \tag{37}$$

Hence, we obtain the time derivatives of the order parameters in a manner similar to that for SGD:

$$\begin{aligned}
\frac{dR_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \left[ \theta_{ik} \langle \delta_k x_j \rangle_z + \langle \delta_k \mathbf{x}^T \rangle_z \Theta_{ik}^T \mathbf{R}_{\cdot j} \right], \\
\frac{dS_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \left[ \theta_{ik} \langle \delta_k x_j \rangle_z + \langle \delta_k \mathbf{x}^T \rangle_z \Theta_{ik}^T \mathbf{Q}_{\cdot j} \right], \\
\frac{d\phi_{ij}}{d\alpha} &= \eta^2 \sum_{k,l=1}^K \theta_{ik} \langle \delta_k \delta_l \rangle_z \theta_{lj}, \\
\frac{dQ_{ij}}{d\alpha} &= \frac{dS_{ij}}{d\alpha} + \frac{dS_{ji}}{d\alpha} + \frac{d\phi_{ij}}{d\alpha},
\end{aligned} \tag{38}$$

where  $\mathbf{R}_{\cdot j}$  denotes the  $j$ th column of  $\mathbf{R}$ , and so on [6].

#### D. Dynamics of $\hat{\mathbf{G}}^{-1}$ for ANG D

In this subsection, we derive the dynamics of the approximated inverse of the Fisher information matrix  $\hat{\mathbf{G}}^{-1}$  (from here on, for simplicity we write  $\mathbf{H}$  instead of  $\hat{\mathbf{G}}^{-1}$ ) in the simplified version of ANG D. Unlike NGD, there are three types of difficulty in ANG D. 1) The dependence between  $\mathbf{J}$  and  $\mathbf{H}$ , because ANG D has two dynamics of  $\mathbf{J}$  and  $\mathbf{H}$  and they share each input  $\boldsymbol{\xi}$ . We introduce an approximation (*two  $\boldsymbol{\xi}$  rule*) to negate this dependence. 2) Higher-order self-correlations of  $\mathbf{H}$ , which originate from the update rule of  $\mathbf{H}$  itself. We negate these infinite correlations by exploiting the  $\rho \ll 1$  assumption of the simplified ANG D, because  $n$ -order self-correlation is scaled by  $O(\rho^n)$ . 3) The high complexity of  $\mathbf{H}$ ;  $\mathbf{G}$  and  $\mathbf{G}^{-1}$  are sufficiently characterized by the subspace of the student weight vectors, whereas  $\mathbf{H}$  is not. This complexity will be managed in the next subsection.

First, we introduce *two  $\boldsymbol{\xi}$  rule*. From Eqs. (8) and (11), the value of  $\mathbf{J}$  and  $\mathbf{H}$  for an infinitesimal time period ( $d\alpha$ ) after current time  $\alpha$  are

$$\begin{aligned}
\mathbf{J}^{(\alpha+d\alpha)} &= \mathbf{J}^{(\alpha)} + \sum_{\mu=0}^{Nd\alpha-1} \Delta \mathbf{J}^{(\tau)} \\
&= \mathbf{J}^{(\alpha)} - \frac{\eta}{N} \sum_{\mu=0}^{Nd\alpha-1} \check{\boldsymbol{\xi}}_{\mathbf{J}^{(\tau)}}(\boldsymbol{\xi}^{(\tau)}) \mathbf{H}^{(\tau)} \nabla f_{(\tau)}^{(\tau)},
\end{aligned} \tag{39}$$

$$\begin{aligned}
\mathbf{H}^{(\alpha+d\alpha)} &= \mathbf{H}^{(\alpha)} + \sum_{\mu=0}^{Nd\alpha-1} \Delta \mathbf{H}^{(\tau)} \\
&= \mathbf{H}^{(\alpha)} + \sum_{\mu=0}^{Nd\alpha-1} \rho' \left[ \mathbf{H}^{(\tau)} - \mathbf{H}^{(\tau)} [\nabla f_{(\tau)}^{(\tau)}] [\nabla f_{(\tau)}^{(\tau)}]^T \mathbf{H}^{(\tau)} \right] \\
&= \{1 + \rho'\}^{Nd\alpha} \mathbf{H}^{(\alpha)} - \sum_{\mu=0}^{Nd\alpha-1} \{1 + \rho'\}^{Nd\alpha-1-\mu} \rho' \mathbf{H}^{(\tau)} [\nabla f_{(\tau)}^{(\tau)}] [\nabla f_{(\tau)}^{(\tau)}]^T \mathbf{H}^{(\tau)},
\end{aligned} \tag{40}$$

where  $\eta$  and  $\rho$  are  $O(1)$  with respect to  $N$ ,  $\rho' \equiv \frac{\rho/N}{1-\rho/N}$ , and  $\tau \equiv \alpha + \mu/N$ , while  $\nabla f_{(\bullet)}^{(\circ)}$  denotes  $\nabla_{\mathbf{J}} f_{\mathbf{J}}(\boldsymbol{\xi}^{(\circ)})|_{\mathbf{J}=\mathbf{J}^{(\circ)}}$ . These two equations show that both  $\mathbf{J}$  and  $\mathbf{H}$  include a common random vector:  $\nabla f_{(\tau)}^{(\tau)}$ . Therefore,  $\mathbf{J}$  and  $\mathbf{H}$  become dependent on each other. To negate this dependence, we introduce a new update rule – we draw two  $\boldsymbol{\xi}$  independently, one for the  $\mathbf{J}$  update and the other for the  $\mathbf{H}$  update in each learning step (*two  $\boldsymbol{\xi}$  rule*) – so that this dependence disappears. Under this *two  $\boldsymbol{\xi}$  rule*, we can fix  $\mathbf{J}$  during  $d\alpha$ , and reduce Eq. (40) to

$$\mathbf{H}^{(\alpha+d\alpha)} = \{1 + \rho'\}^{Nd\alpha} \mathbf{H}^{(\alpha)}$$

$$- \sum_{\mu=0}^{Nd\alpha-1} \{1 + \rho'\}^{Nd\alpha-1-\mu} \mathbf{H}^{(\tau)} \mathbf{F}^{(\tau)} \mathbf{H}^{(\tau)}, \tag{41}$$

where  $\mathbf{F}^{(\tau)} \equiv \rho' [\nabla f_{(\alpha)}^{(\tau)}] [\nabla f_{(\alpha)}^{(\tau)}]^T$ . This *two  $\boldsymbol{\xi}$  rule* will be validated in Section IV.

Next, we negate the higher-order self-correlations of  $\mathbf{H}$ . We notice that this Eq. (41) is still difficult to solve, because it includes highly self-correlated terms with respect to the old random matrix  $\mathbf{F}^{(\tau)}$ , ( $\alpha \leq \tau' < \tau$ ). For example, the most

correlated term is

$$\mathbf{H}^{(\alpha)} \mathbf{F}^{(\alpha)} \mathbf{H}^{(\alpha)} \mathbf{F}^{(\alpha+\frac{1}{N})} \mathbf{H}^{(\alpha)} \mathbf{F}^{(\alpha)} \mathbf{H}^{(\alpha)} \mathbf{F}^{(\alpha+\frac{2}{N})} \dots \quad (42)$$

However, we can ignore these terms including  $n$ -order self-correlation ( $n \geq 2$ ), because they are at most  $O(\rho^n)$  and small enough under the assumption of  $\rho \ll 1$  for the simplified ANG. Therefore, we sum up only the  $O(\rho^0)$  and  $O(\rho^1)$  terms:

$$\begin{aligned} \mathbf{H}^{(\alpha+d\alpha)} &= \{1+\rho'\}^{Nd\alpha} \mathbf{H}^{(\alpha)} + \sum_{\mu=0}^{Nd\alpha-1} \{1+\rho'\}^{Nd\alpha-1-\mu} \\ &\quad \times \left[ \{1+\rho'\}^\mu \mathbf{H}^{(\alpha)} \right] \mathbf{F}^{(\tau)} \left[ \{1+\rho'\}^\mu \mathbf{H}^{(\alpha)} \right]. \end{aligned} \quad (43)$$

At the large limit of  $N$ , Eq. (43) becomes

$$\mathbf{H}^{(\alpha+d\alpha)} = e^{\rho d\alpha} \left[ \mathbf{H}^{(\alpha)} + \{1-e^{\rho d\alpha}\} \mathbf{H}^{(\alpha)} \mathbf{G} \mathbf{H}^{(\alpha)} \right], \quad (44)$$

where we use  $\sum \mathbf{F} \rightarrow \mathbf{G}$  (see Appendix A for further discussion about the convergence of  $\mathbf{F}$ ). Then, we obtain the dynamics of  $\mathbf{H}$ :

$$\begin{aligned} \frac{d\mathbf{H}}{d\alpha} &= \lim_{d\alpha \rightarrow 0} \frac{\mathbf{H}^{(\alpha+d\alpha)} - \mathbf{H}^{(\alpha)}}{d\alpha} \\ &= \rho [\mathbf{H} - \mathbf{H} \mathbf{G} \mathbf{H}]. \end{aligned} \quad (45)$$

We also obtain the usual order parameter dynamics. As the update equation of each  $\mathbf{J}_i$  for ANG is

$$\Delta \mathbf{J}_i = -\frac{\eta}{N} \sum_{k=1}^K \delta_k \mathbf{H}_{ik} \boldsymbol{\xi}, \quad (46)$$

we can easily get

$$\begin{aligned} \frac{dR_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \langle \delta_k \boldsymbol{\xi}^T \mathbf{H}_{ik}^T \mathbf{B}_j \rangle_{\boldsymbol{\xi}}, \\ \frac{dS_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \langle \delta_k \boldsymbol{\xi}^T \mathbf{H}_{ik}^T \mathbf{J}_j \rangle_{\boldsymbol{\xi}}, \\ \frac{d\phi_{ij}}{d\alpha} &= \frac{\eta^2}{N} \sum_{k,l=1}^K \langle \delta_k \delta_l \boldsymbol{\xi}^T \mathbf{H}_{ik}^T \mathbf{H}_{jl} \boldsymbol{\xi} \rangle_{\boldsymbol{\xi}}, \\ \frac{dQ_{ij}}{d\alpha} &= \frac{dS_{ij}}{d\alpha} + \frac{dS_{ji}}{d\alpha} + \frac{d\phi_{ij}}{d\alpha}. \end{aligned} \quad (47)$$

These order parameter dynamics still include  $N$ -dimensional vectors and  $N \times N$  matrices, but we obtain the order parameter dynamics expressed by the order parameters themselves in the next subsection.

### E. Order parameter representation for $\mathbf{H} \equiv \hat{\mathbf{G}}^{-1}$

In this subsection, we extract the new order parameters from  $\mathbf{H}$ . To characterize  $\mathbf{H}$ , we should consider the history of  $\mathbf{J}$ , which includes two types of movement. 1) Approach

to the corresponding teacher vectors. To deal with this movement, we use not only the student weight vectors but also the teacher weight vectors to express  $\mathbf{H}$ , although  $\mathbf{G}$  and  $\mathbf{G}^{-1}$  are sufficiently expressed by only the current student weight vectors. 2) Escape from the subspace made by all the weight vectors. The direction of this movement is random, although its speed is deterministic as specified by the  $d\phi$  term. The order parameter expressions of SGD and NGD do not suffer from this randomness because they discard it by exploiting the rotation invariance of the system. In ANG, however, this randomness upsets the order parameter expression of  $\mathbf{H}$ ; it produces a hysteresis component of  $\mathbf{H}$  or the residual fluctuating term ( $\mathbf{E}$  as introduced below) when  $\mathbf{H}$  is expressed through order parameters. This fluctuation term is negligible at the large limit of  $N$ , but its square is not. Here, we realistically consider the  $\mathbf{H}$  dynamics by taking the second powers of this fluctuation term into account, and obtain an effective order parameter expression of  $\mathbf{H}$ . This theoretical result will be numerically validated in the next section.

First, we express  $\mathbf{H}$  in a similar manner to  $\mathbf{G}$  in Eq. (27) as

$$\mathbf{H}_{ij} = \omega_{ij} \mathbf{I} + [\mathbf{U}, \mathbf{V}] \begin{bmatrix} \boldsymbol{\Omega}_{ij} & \boldsymbol{\Omega}'_{ij} \\ \boldsymbol{\Omega}''_{ij} & \boldsymbol{\Omega}'''_{ij} \end{bmatrix} [\mathbf{U}, \mathbf{V}]^T \quad (48)$$

$$= \omega_{ij} \mathbf{I} + \mathbf{U} \boldsymbol{\Omega}_{ij} \mathbf{U}^T + \mathbf{E}_{ij}, \quad (49)$$

$$\mathbf{E}_{ij} \equiv [\mathbf{U}, \mathbf{V}] \begin{bmatrix} \mathbf{0} & \boldsymbol{\Omega}'_{ij} \\ \boldsymbol{\Omega}''_{ij} & \boldsymbol{\Omega}'''_{ij} \end{bmatrix} [\mathbf{U}, \mathbf{V}]^T, \quad (50)$$

and prove  $\mathbf{E}_{ij}$  is negligible. Note that we use not only the student weight vectors, but also the teacher weight vectors; i.e., we use  $\mathbf{U} \equiv [\mathbf{J}_1, \dots, \mathbf{J}_K, \mathbf{B}_1, \dots, \mathbf{B}_M] \in \mathbb{R}^{N \times (K+M)}$  instead of  $\mathbf{J}' \equiv [\mathbf{J}_1, \dots, \mathbf{J}_K]$ . This is because the student vectors move toward the corresponding teacher vectors, and  $\mathbf{H}$  holds the component made from the old student vectors. Here,  $\mathbf{V} \in \mathbb{R}^{N \times (N-K-M)}$  is re-defined as the orthonormal bases of the null space of  $\mathbf{U}$  rather than  $\mathbf{J}'$ , which satisfies  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$  and  $\mathbf{U}^T \mathbf{V} = \mathbf{0}$ . Here,  $\omega_{ij} \in \mathbb{R}$  and  $\boldsymbol{\Omega}_{ij} \in \mathbb{R}^{(K+M) \times (K+M)}$  are the candidates for the new order parameters. Also,  $\boldsymbol{\Omega}'_{ij} \in \mathbb{R}^{(K+M) \times (N-K-M)}$ ,  $\boldsymbol{\Omega}''_{ij} \in \mathbb{R}^{(N-K-M) \times (K+M)}$ , and  $\boldsymbol{\Omega}'''_{ij} \in \mathbb{R}^{(N-K-M) \times (N-K-M)}$  are large matrices. As  $\mathbf{H} = [\mathbf{H}_{ij}]_{i,j=1,\dots,K}$  is a symmetric matrix,  $\boldsymbol{\omega} \equiv [\omega_{ij}]_{i,j=1,\dots,K}$ ,  $\boldsymbol{\Omega} \equiv [\boldsymbol{\Omega}_{ij}]_{i,j=1,\dots,K}$ , and  $\boldsymbol{\Omega}''' \equiv [\boldsymbol{\Omega}'''_{ij}]_{i,j=1,\dots,K}$  are also symmetric matrices, while  $\boldsymbol{\Omega}' \equiv [\boldsymbol{\Omega}'_{ij}]_{i,j=1,\dots,K}$  and  $\boldsymbol{\Omega}'' \equiv [\boldsymbol{\Omega}''_{ij}]_{i,j=1,\dots,K}$  are symmetric with respect to each other.

Next, we find appropriate dynamics of  $\boldsymbol{\omega}$ ,  $\boldsymbol{\Omega}$ ,  $\boldsymbol{\Omega}'$ ,  $\boldsymbol{\Omega}''$ , and  $\boldsymbol{\Omega}'''$  that satisfy the dynamics of  $\mathbf{H}$  given by Eq. (45). For convenience, we consider an infinitesimal change of  $\mathbf{H}_{ij}$  from Eq. (45),

$$d\mathbf{H}_{ij} = \rho [\mathbf{H} - \mathbf{H} \mathbf{G} \mathbf{H}]_{ij} d\alpha \quad (51)$$

By substituting Eq. (48) for  $\mathbf{H}$ , we easily obtain the decomposed form of  $d\mathbf{H}_{ij}$  as

$$d\mathbf{H}_{ij} = \gamma_{ij} \mathbf{I} + [\mathbf{U}, \mathbf{V}] \begin{bmatrix} \boldsymbol{\Gamma}_{ij} & \boldsymbol{\Gamma}'_{ij} \\ \boldsymbol{\Gamma}''_{ij} & \boldsymbol{\Gamma}'''_{ij} \end{bmatrix} [\mathbf{U}, \mathbf{V}]^T, \quad (52)$$

where

$$\gamma_{ij} = \rho[\omega - \omega\lambda\omega]_{ij}d\alpha, \quad (53)$$

$$\Gamma_{ij} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega} - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega} + \tilde{\omega}\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} d\alpha, \quad (54)$$

$$\Gamma'_{ij} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega}' - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega}' - \tilde{\Omega}'\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} d\alpha, \quad (55)$$

$$\Gamma''_{ij} = \rho \left[ \tilde{\Omega}'' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega}\tilde{\lambda}\tilde{\Omega}'' - \tilde{\Omega}''\tilde{\lambda}\tilde{\Omega}''' \right]_{ij} C^{-\frac{1}{2}} d\alpha, \quad (56)$$

$$\Gamma'''_{ij} = \rho \left[ \tilde{\Omega}''' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega}' - \tilde{\omega}\tilde{\lambda}\tilde{\Omega}''' - \tilde{\Omega}''\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}''\tilde{\lambda}\tilde{\Omega}''' \right]_{ij} d\alpha, \quad (57)$$

where  $\tilde{\Omega} \equiv [\tilde{\Omega}_{ij} \equiv C^{\frac{1}{2}}\Omega_{ij}C^{\frac{1}{2}} + \omega_{ij}\mathbf{I}]_{i,j=1,\dots,K}$ ,  $\tilde{\Omega}' \equiv [\tilde{\Omega}'_{ij} \equiv C^{\frac{1}{2}}\Omega'_{ij}]_{i,j=1,\dots,K}$ , and  $\tilde{\Omega}'' \equiv [\tilde{\Omega}''_{ij} \equiv \Omega_{ij}C^{\frac{1}{2}}]_{i,j=1,\dots,K}$  are introduced to simplify multiplications, while,  $\tilde{\lambda} = [\lambda_{ij}\mathbf{I}]_{i,j=1,\dots,K}$  and  $\tilde{\omega} = [\omega_{ij}\mathbf{I}]_{i,j=1,\dots,K}$  are matrices extended to an appropriate size; e.g., the same size as that of  $\Omega$  or  $\Omega'''$ . The  $\tilde{\Lambda}$  is re-defined as  $\tilde{\Lambda} \equiv [\tilde{\Lambda}_{ij} \equiv C^{\frac{1}{2}} \begin{bmatrix} \Lambda_{ij} & 0 \\ 0 & 0 \end{bmatrix} C^{\frac{1}{2}} + \lambda_{ij}\mathbf{I}]_{i,j=1,\dots,K}$ .

The infinitesimal change of Eq. (48) is given by

$$d\mathbf{H}_{ij} = \{d\omega_{ij}\}\mathbf{I} + \mathbf{W} \begin{bmatrix} \Omega_{ij} + d\Omega_{ij} & \Omega'_{ij} + d\Omega'_{ij} \\ \Omega''_{ij} + d\Omega''_{ij} & \Omega'''_{ij} + d\Omega'''_{ij} \end{bmatrix} \mathbf{W}^T - [\mathbf{U}, \mathbf{V}] \begin{bmatrix} \Omega_{ij} & \Omega'_{ij} \\ \Omega''_{ij} & \Omega'''_{ij} \end{bmatrix} [\mathbf{U}, \mathbf{V}]^T, \quad (58)$$

where  $\mathbf{W} \equiv [\mathbf{U} + d\mathbf{U}, \mathbf{V} + d\mathbf{V}] \in \mathbb{R}^{N \times N}$ . Note that the conventional total differential is not applicable because  $\mathbf{U}$  is not differentiable. We also decompose this matrix by  $\mathbf{W}$  rather than  $[\mathbf{U}, \mathbf{V}]$  because  $\mathbf{U}$  and  $\mathbf{V}$  are also moving,

$$d\mathbf{H}_{ij} = l_{ij}\mathbf{I} + \mathbf{W} \begin{bmatrix} L_{ij} & L'_{ij} \\ L''_{ij} & L'''_{ij} \end{bmatrix} \mathbf{W}^T. \quad (59)$$

Here,  $L_{ij}, L'_{ij}, L''_{ij}$ , and  $L'''_{ij}$  are easily determined by calculating

$$\begin{bmatrix} L_{ij} & L'_{ij} \\ L''_{ij} & L'''_{ij} \end{bmatrix} \equiv [\mathbf{W}^T \mathbf{W}]^{-1} \mathbf{W}^T \left[ d\mathbf{H}_{ij} - \{d\omega_{ij}\}\mathbf{I} \right] \times \mathbf{W} [\mathbf{W}^T \mathbf{W}]^{-1}. \quad (60)$$

Note that  $\mathbf{W}[\mathbf{W}^T \mathbf{W}]^{-1} \mathbf{W}^T$  is an identity matrix. By erasing small terms such as  $O(d^2)$ , we obtain the following equation,

$$d\mathbf{H}_{ij} = l_{ij}\mathbf{I} + [\mathbf{U}, \mathbf{V}] \begin{bmatrix} L_{ij} & L'_{ij} \\ L''_{ij} & L'''_{ij} \end{bmatrix} [\mathbf{U}, \mathbf{V}]^T, \quad (61)$$

where

$$l_{ij} = d\omega_{ij}, \quad (62)$$

$$L_{ij} = d\Omega_{ij} + C^{-1} [dC - d\mathbf{O}] \Omega_{ij} + \Omega_{ij} [dC - d\mathbf{O}]^T C^{-1} - C^{-1} [d\mathbf{U}]^T \mathbf{V} \Omega''_{ij} - \Omega'_{ij} \mathbf{V}^T [d\mathbf{U}] C^{-1}, \quad (63)$$

$$L'_{ij} = d\Omega'_{ij} + C^{-1} [dC - d\mathbf{O}] \Omega'_{ij} - \Omega'_{ij} \mathbf{V}^T [d\mathbf{V}] - C^{-1} [d\mathbf{U}]^T \mathbf{V} \Omega'''_{ij} - \Omega_{ij} \mathbf{U}^T [d\mathbf{V}], \quad (64)$$

$$L''_{ij} = d\Omega''_{ij} - [d\mathbf{V}]^T \mathbf{V} \Omega''_{ij} + \Omega''_{ij} [dC - d\mathbf{O}]^T C^{-1} - [d\mathbf{V}]^T \mathbf{U} \Omega_{ij} - \Omega'''_{ij} \mathbf{V}^T [d\mathbf{U}] C^{-1}, \quad (65)$$

$$L'''_{ij} = d\Omega'''_{ij} - [d\mathbf{V}]^T \mathbf{V} \Omega'''_{ij} - \Omega'''_{ij} \mathbf{V}^T [d\mathbf{V}] - [d\mathbf{V}]^T \mathbf{U} \Omega'_{ij} - \Omega''_{ij} \mathbf{U}^T [d\mathbf{V}], \quad (66)$$

where we used

$$\mathbf{W}^T \mathbf{W} = \begin{bmatrix} C + dC & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad (67)$$

$$\mathbf{W}^T [\mathbf{U}, \mathbf{V}] = \begin{bmatrix} C + dC & [d\mathbf{U}]^T \mathbf{V} \\ [d\mathbf{V}]^T \mathbf{U} & \mathbf{I} + [d\mathbf{V}]^T \mathbf{V} \end{bmatrix}, \quad (68)$$

where  $C \equiv \mathbf{U}^T \mathbf{U} = \begin{bmatrix} Q & R \\ R^T & T \end{bmatrix} \in \mathbb{R}^{(K+M) \times (K+M)}$ ,  $dC \equiv [\mathbf{U} + d\mathbf{U}]^T [\mathbf{U} +$

$$d\mathbf{U}] - C = \begin{bmatrix} dQ & dR \\ dR^T & dT \end{bmatrix}, \text{ and } d\mathbf{O} \equiv [d\mathbf{U}]^T \mathbf{U} = \begin{bmatrix} dS & dR \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

To realize  $d\mathbf{H}$  of Eq. (52) through Eq. (61) under the given  $\mathbf{U}, \mathbf{V}, d\mathbf{U}$ , and  $d\mathbf{V}$ ,  $l_{ij} = \gamma_{ij}$  and  $\begin{bmatrix} L_{ij} & L'_{ij} \\ L''_{ij} & L'''_{ij} \end{bmatrix} = \begin{bmatrix} \Gamma_{ij} & \Gamma'_{ij} \\ \Gamma''_{ij} & \Gamma'''_{ij} \end{bmatrix}$  should be satisfied. We then easily obtain an appropriate infinitesimal change of  $\omega, \Omega, \Omega', \Omega''$ , and  $\Omega'''$  as,

$$d\omega_{ij} = \rho[\omega - \omega\lambda\omega]_{ij}d\alpha, \quad (69)$$

$$d\Omega_{ij} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega} - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega} + \tilde{\omega}\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} d\alpha \\ - C^{-1} [dC - dO] \Omega_{ij} - \Omega_{ij} [dC - dO]^T C^{-1} + C^{-1} [dU]^T V \Omega_{ij}' + \Omega_{ij}' V^T [dU] C^{-1}, \quad (70)$$

$$d\Omega'_{ij} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega}' - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega}' - \tilde{\Omega}'\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} d\alpha \\ - C^{-1} [dC - dO] \Omega'_{ij} + \Omega'_{ij} V^T [dV] + C^{-1} [dU]^T V \Omega_{ij}'' + \Omega_{ij}'' V^T [dV], \quad (71)$$

$$d\Omega''_{ij} = \rho \left[ \tilde{\Omega}'' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega}\tilde{\lambda}\tilde{\Omega}'' - \Omega''\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} d\alpha \\ + [dV]^T V \Omega_{ij}'' - \Omega_{ij}'' [dC - dO]^T C^{-1} + [dV]^T U \Omega_{ij} + \Omega_{ij}'' V^T [dU] C^{-1}, \quad (72)$$

$$d\Omega'''_{ij} = \rho \left[ \Omega''' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega}' - \tilde{\omega}\tilde{\lambda}\Omega''' - \Omega'''\tilde{\lambda}\tilde{\omega} - \Omega'''\tilde{\lambda}\Omega'' \right]_{ij} d\alpha \\ + [dV]^T V \Omega_{ij}''' + \Omega_{ij}''' V^T [dV] + [dV]^T U \Omega'_{ij} + \Omega'_{ij} U^T [dV]. \quad (73)$$

We notice that, with respect to  $N$ , the orders of each element of the matrices  $\omega, \Omega, \Omega', \Omega''$ , and  $\Omega'''$  are at most  $O(1), O(1), O(\frac{1}{\sqrt{N}}), O(\frac{1}{\sqrt{N}})$ , and  $O(\frac{1}{N})$  (see Appendix B). This suggests that  $\mathbf{E}$  is negligible at the large limit of  $N$ . Actually, if we explicitly assume that the escape directions of the student

vectors are completely random,  $\langle [I - UC^{-1}U^T]dU \rangle_{\xi} = \mathbf{0}$  or equivalently  $\langle V^T dU \rangle_{\xi} = \mathbf{0}$  or  $\langle dV \rangle_{\xi} = \mathbf{0}$ , the time integrations of the terms including  $U^T dV, V^T dU$ , and  $V^T dV$  converge to zero, and we obtain

$$\frac{d\omega_{ij}}{d\alpha} = \rho[\omega - \omega\lambda\omega]_{ij}, \quad (74)$$

$$\frac{d\Omega_{ij}}{d\alpha} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega} - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega} + \tilde{\omega}\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} - C^{-1} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right]^T \Omega_{ij} - \Omega_{ij} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right] C^{-1}, \quad (75)$$

$$\frac{d\Omega'_{ij}}{d\alpha} = \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega}' - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega}' - \tilde{\Omega}'\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} - C^{-1} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right]^T \Omega'_{ij}, \quad (76)$$

$$\frac{d\Omega''_{ij}}{d\alpha} = \rho \left[ \tilde{\Omega}'' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega}\tilde{\lambda}\tilde{\Omega}'' - \Omega''\tilde{\lambda}\tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} - \Omega_{ij}'' \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right] C^{-1}, \quad (77)$$

$$\frac{d\Omega'''_{ij}}{d\alpha} = \rho \left[ \Omega''' - \tilde{\Omega}''\tilde{\Lambda}\tilde{\Omega}' - \tilde{\omega}\tilde{\lambda}\Omega''' - \Omega'''\tilde{\lambda}\tilde{\omega} - \Omega'''\tilde{\lambda}\Omega'' \right]_{ij}. \quad (78)$$

These equations guarantee that  $\Omega', \Omega'',$  and  $\Omega'''$  are always zero at the large limit of  $N$ , because the initial zero values of  $\Omega'_{ij}, \Omega''_{ij},$  and  $\Omega'''_{ij}$  are preserved. Consequently, we adopt  $\omega$  and  $\Omega$  as the new order parameters, which successfully express  $\mathbf{H}$ .

Next, we consider the  $HGH$  term in Eq. (51) because a small fluctuation of  $\mathbf{E}$  might become significant in this second-order of the  $\mathbf{H}$  term. Specifically, we evaluate the square of the fluctuation term  $\mathbf{E}$  (to be exact,  $\omega$  and  $\Omega$  also fluctuate, but their fluctuations do not become significant because these matrices are sufficiently small compared with  $N$ ).

In a similar way as above, we let

$$[E^2]_{ij} = [U, V] \begin{bmatrix} \Upsilon_{ij} & \Upsilon'_{ij} \\ \Upsilon''_{ij} & \Upsilon'''_{ij} \end{bmatrix} [U, V]^T, \quad (79)$$

where

$$\begin{bmatrix} \Upsilon_{ij} & \Upsilon'_{ij} \\ \Upsilon''_{ij} & \Upsilon'''_{ij} \end{bmatrix} \equiv \sum_{k=1}^K \left[ \begin{matrix} \Omega'_{ik} \Omega''_{kj} & \Omega'_{ik} \Omega'''_{kj} \\ \Omega''_{ik} \Omega'_{kj} & \Omega''_{ik} C \Omega'_{kj} + \Omega'''_{ik} \Omega''_{kj} \end{matrix} \right]. \quad (80)$$

Then, an infinitesimal change of  $E^2$  is

$$d[E^2]_{ij} = [E + dE]^2 - E^2 \Big|_{ij}$$



$$\begin{aligned}
&= \sum_{k=1}^K \mathbf{W} \begin{bmatrix} [\boldsymbol{\Omega}'_{ik} + d\boldsymbol{\Omega}'_{ik}][\boldsymbol{\Omega}''_{kj} + d\boldsymbol{\Omega}''_{kj}] & [\boldsymbol{\Omega}'_{ik} + d\boldsymbol{\Omega}'_{ik}][\boldsymbol{\Omega}'''_{kj} + d\boldsymbol{\Omega}'''_{kj}] \\ [\boldsymbol{\Omega}'''_{ik} + d\boldsymbol{\Omega}'''_{ik}][\boldsymbol{\Omega}''_{kj} + d\boldsymbol{\Omega}''_{kj}] & [\boldsymbol{\Omega}''_{ik} + d\boldsymbol{\Omega}''_{ik}][C + dC][\boldsymbol{\Omega}'_{kj} + d\boldsymbol{\Omega}'_{kj}] + [\boldsymbol{\Omega}'''_{ik} + d\boldsymbol{\Omega}'''_{ik}][\boldsymbol{\Omega}'''_{kj} + d\boldsymbol{\Omega}'''_{kj}] \end{bmatrix} \mathbf{W}^T \\
&\quad - \sum_{k=1}^K [U, V] \begin{bmatrix} \boldsymbol{\Omega}'_{ik} \boldsymbol{\Omega}''_{kj} & \boldsymbol{\Omega}'_{ik} \boldsymbol{\Omega}'''_{kj} \\ \boldsymbol{\Omega}'''_{ik} \boldsymbol{\Omega}''_{kj} & \boldsymbol{\Omega}''_{ik} C \boldsymbol{\Omega}'_{kj} + \boldsymbol{\Omega}'''_{ik} \boldsymbol{\Omega}'''_{kj} \end{bmatrix} [U, V]^T. \tag{81}
\end{aligned}$$

Substituting Eqs. (70)-(73) for  $d\boldsymbol{\Omega}'$ ,  $d\boldsymbol{\Omega}''$ , and  $d\boldsymbol{\Omega}'''$ , and erasing small terms, we obtain an appropriate infinitesimal change of  $\boldsymbol{\Upsilon}_{ij}$ ,  $\boldsymbol{\Upsilon}'_{ij}$ ,  $\boldsymbol{\Upsilon}''_{ij}$ , and  $\boldsymbol{\Upsilon}'''_{ij}$ . Briefly, the point is that  $[d\boldsymbol{\Omega}'_{ik}][d\boldsymbol{\Omega}''_{kj}]$  is not negligible but  $\boldsymbol{\Omega}_{ik} \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix} \boldsymbol{\Omega}_{kj}$ , because  $U^T [dV][dV]^T U = \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix}$  (see Appendix B). Actually, we obtain,

$$\begin{aligned}
d\boldsymbol{\Upsilon}_{ij} &= \rho C^{-\frac{1}{2}} \left[ 2\tilde{\boldsymbol{\Omega}}' \tilde{\boldsymbol{\Omega}}'' - \tilde{\boldsymbol{\Omega}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Omega}}' \tilde{\boldsymbol{\Omega}}'' - \tilde{\boldsymbol{\Omega}}' \tilde{\boldsymbol{\Omega}}'' \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Omega}} \right. \\
&\quad \left. - \tilde{\boldsymbol{\Omega}}' [\tilde{\boldsymbol{\lambda}} \tilde{\omega} + \tilde{\omega} \tilde{\boldsymbol{\lambda}}] \tilde{\boldsymbol{\Omega}}'' \right]_{ij} C^{-\frac{1}{2}} d\alpha \\
&\quad + \sum_{k=1}^K \left[ \boldsymbol{\Gamma}'_{ik} \boldsymbol{\Omega}''_{kj} - C^{-1} [dC - d\boldsymbol{O}] \boldsymbol{\Omega}'_{ik} \boldsymbol{\Omega}''_{kj} + \boldsymbol{\Omega}'_{ik} \boldsymbol{\Gamma}''_{kj} \right. \\
&\quad \left. - \boldsymbol{\Omega}'_{ik} \boldsymbol{\Omega}''_{kj} [dC - d\boldsymbol{O}]^T C^{-1} + \boldsymbol{\Omega}_{ik} \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix} \boldsymbol{\Omega}_{kj} \right]. \tag{82}
\end{aligned}$$

We notice that each element of  $\boldsymbol{\Upsilon}_{ij}$  is initially zero, and becomes  $O(1)$ , because  $\boldsymbol{\Omega}_{ik} \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix} \boldsymbol{\Omega}_{kj}$  is  $O(1)$ , whereas the other terms,  $\boldsymbol{\Upsilon}'_{ij}$ ,  $\boldsymbol{\Upsilon}''_{ij}$ , and  $\boldsymbol{\Upsilon}'''_{ij}$ , are considered to be always zero at the large limit of  $N$ .

We then let the following positive definite and symmetric matrix

$$\boldsymbol{\Upsilon} = \boldsymbol{\Omega}' \boldsymbol{\Omega}'' \in \mathbb{R}^{K(K+M) \times K(K+M)}, \tag{83}$$

be another new order parameter, where  $\boldsymbol{\Upsilon} = [\boldsymbol{\Upsilon}_{ij}]_{i,j=1,\dots,K}$ . In short, we can write

$$\boldsymbol{E}_{ij} \rightarrow \mathbf{0}, \tag{84}$$

$$[\boldsymbol{E}^2]_{ij} - U \boldsymbol{\Upsilon}_{ij} U^T \rightarrow \mathbf{0}, \tag{85}$$

or, equivalently,

$$\omega_{ij} \mathbf{I} + U \boldsymbol{\Omega}_{ij} U^T \rightarrow \boldsymbol{H}_{ij}, \tag{86}$$

$$\begin{aligned}
&\sum_{k=1}^K [\omega_{ik} \mathbf{I} + U \boldsymbol{\Omega}_{ik} U^T][\omega_{kj} \mathbf{I} + U \boldsymbol{\Omega}_{kj} U^T] + U \boldsymbol{\Upsilon}_{ij} U^T \\
&\quad \rightarrow [\boldsymbol{H}^2]_{ij}, \tag{87}
\end{aligned}$$

at the large limit of  $N$ .

This  $\boldsymbol{\Upsilon}$  is useful for estimating the matrices  $\boldsymbol{\Omega}'$  and  $\boldsymbol{\Omega}''$ , which are required for both the  $\boldsymbol{\Omega}'[\tilde{\boldsymbol{\lambda}}\tilde{\omega} + \tilde{\omega}\tilde{\boldsymbol{\lambda}}]\boldsymbol{\Omega}''$  term in Eq. (82) and the  $\tilde{\boldsymbol{\Omega}}'\tilde{\boldsymbol{\lambda}}\tilde{\boldsymbol{\Omega}}''$  term in Eq. (75). Intrinsically, an exact restoration of  $\boldsymbol{\Omega}'$  and  $\boldsymbol{\Omega}''$  using  $\boldsymbol{\Upsilon}$  is impossible, because the number of elements of  $\boldsymbol{\Omega}'$  and  $\boldsymbol{\Omega}''$  is  $O(N)$ , while that of  $\boldsymbol{\Upsilon}$  is  $O(1)$ . Hence, we substitute one of the probable candidates of  $\boldsymbol{\Omega}'_{ij}$  and  $\boldsymbol{\Omega}''_{ij}$  for the true  $\boldsymbol{\Omega}'$  and  $\boldsymbol{\Omega}''$ ; i.e., we use

$$\boldsymbol{\Omega}'_{ij} = \boldsymbol{\Omega}''_{ji}{}^T = \frac{1}{\sqrt{n}} \left[ [\sqrt{\boldsymbol{\Upsilon}}]_{ij}, \dots, [\sqrt{\boldsymbol{\Upsilon}}]_{ij} \right], \tag{88}$$

where  $n \equiv \frac{N-K-M}{K+M}$  was assumed to be a natural number. This Eq. (88) certainly satisfies the necessary condition given by Eq. (83). Substituting this Eq. (88) for Eq. (82), we obtain

$$\begin{aligned}
d\boldsymbol{\Upsilon}_{ij} &= \rho \left[ C^{-\frac{1}{2}} \left[ 2\tilde{\boldsymbol{\Upsilon}} - \tilde{\boldsymbol{\Omega}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Upsilon}} - \tilde{\boldsymbol{\Upsilon}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Omega}} \right]_{ij} C^{-\frac{1}{2}} \right. \\
&\quad \left. - \left[ \sqrt{\boldsymbol{\Upsilon}} [\tilde{\boldsymbol{\lambda}} \tilde{\omega} + \tilde{\omega} \tilde{\boldsymbol{\lambda}}] \sqrt{\boldsymbol{\Upsilon}} \right]_{ij} \right] d\alpha \\
&\quad - C^{-1} [dC - d\boldsymbol{O}] \boldsymbol{\Upsilon}_{ij} - \boldsymbol{\Upsilon}_{ij} [dC - d\boldsymbol{O}]^T C^{-1} \\
&\quad + \sum_{k=1}^K \boldsymbol{\Omega}_{ik} \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix} \boldsymbol{\Omega}_{kj}, \tag{89}
\end{aligned}$$

where  $\tilde{\boldsymbol{\Upsilon}} \equiv [\tilde{\boldsymbol{\Upsilon}}_{ij} \equiv C^{\frac{1}{2}} \boldsymbol{\Upsilon}_{ij} C^{\frac{1}{2}}]_{i,j=1,\dots,K}$ . In a similar manner, we can substitute Eq. (88) for Eq. (75).

Finally, we obtain the new order parameter dynamics through Eqs. (74), (75), and (89) as

$$\frac{d\omega_{ij}}{d\alpha} = \rho [\omega - \omega \boldsymbol{\lambda} \omega]_{ij}, \tag{90}$$

$$\begin{aligned}
\frac{d\boldsymbol{\Omega}_{ij}}{d\alpha} &= \rho \left[ C^{-\frac{1}{2}} \left[ \tilde{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Omega}} - \tilde{\omega} + \tilde{\omega} \tilde{\boldsymbol{\lambda}} \tilde{\omega} \right]_{ij} C^{-\frac{1}{2}} - \left[ \sqrt{\boldsymbol{\Upsilon}} \tilde{\boldsymbol{\lambda}} \sqrt{\boldsymbol{\Upsilon}} \right]_{ij} \right] \\
&\quad - C^{-1} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right]^T \boldsymbol{\Omega}_{ij} - \boldsymbol{\Omega}_{ij} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right] C^{-1}, \tag{91}
\end{aligned}$$

$$\begin{aligned}
\frac{d\boldsymbol{\Upsilon}_{ij}}{d\alpha} &= \rho \left[ C^{-\frac{1}{2}} \left[ 2\tilde{\boldsymbol{\Upsilon}} - \tilde{\boldsymbol{\Omega}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Upsilon}} - \tilde{\boldsymbol{\Upsilon}} \tilde{\boldsymbol{\Lambda}} \tilde{\boldsymbol{\Omega}} \right]_{ij} C^{-\frac{1}{2}} - \left[ \sqrt{\boldsymbol{\Upsilon}} [\tilde{\boldsymbol{\lambda}} \tilde{\omega} + \tilde{\omega} \tilde{\boldsymbol{\lambda}}] \sqrt{\boldsymbol{\Upsilon}} \right]_{ij} \right] \\
&\quad - C^{-1} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right]^T \boldsymbol{\Upsilon}_{ij} - \boldsymbol{\Upsilon}_{ij} \left[ \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} \frac{dR}{d\alpha} \right] C^{-1}
\end{aligned}$$

$$+ \sum_{k=1}^K \mathbf{\Omega}_{ik} \begin{bmatrix} \frac{d\phi}{d\alpha} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{\Omega}_{kj}. \quad (92)$$

Intuitively, in Eqs. (90)-(92), the terms including  $\rho$  correspond to the dynamics of  $\mathbf{H}$ , whereas the terms including  $\begin{bmatrix} \frac{dS}{d\alpha} + \frac{d\phi}{d\alpha} & \frac{dR}{d\alpha} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$  keep the component of  $\mathbf{H}$  constant which is still expressible by  $\mathbf{\Omega}$  or  $\Upsilon$  under the movement of  $\mathbf{U}$ . The term including  $\begin{bmatrix} \frac{d\phi}{d\alpha} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$  means the component of  $\mathbf{H}$  which is no longer expressible by  $\mathbf{\Omega}$  under the escape movement of  $\mathbf{U}$  into the null space of  $\mathbf{U}$ .

The dynamics of the usual order parameters for ANG D given by Eq. (47) can be rewritten using the new order parameters in a manner similar to that for NGD,

$$\begin{aligned} \frac{dR_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \left[ \omega_{ik} \langle \delta_k x_j \rangle_z + \langle \delta_k z^T \rangle_z \mathbf{\Omega}_{ik}^T \begin{bmatrix} \mathbf{R} \\ \mathbf{T} \end{bmatrix}_{\bullet j} \right], \\ \frac{dS_{ij}}{d\alpha} &= -\eta \sum_{k=1}^K \left[ \omega_{ik} \langle \delta_k x_j \rangle_z + \langle \delta_k z^T \rangle_z \mathbf{\Omega}_{ik}^T \begin{bmatrix} \mathbf{Q} \\ \mathbf{R}^T \end{bmatrix}_{\bullet j} \right], \\ \frac{d\phi_{ij}}{d\alpha} &= \eta^2 \sum_{k,l=1}^K \omega_{ik} \langle \delta_k \delta_l \rangle_z \omega_{lj}, \\ \frac{dQ_{ij}}{d\alpha} &= \frac{dS_{ij}}{d\alpha} + \frac{dS_{ji}}{d\alpha} + \frac{d\phi_{ij}}{d\alpha}. \end{aligned} \quad (93)$$

Note that these usual order parameter dynamics are not affected by the small fluctuation,  $\mathbf{E}$  (see Appendix C). Equations (90)-(93) are the order parameter dynamics expressed by the order parameters themselves.

#### IV. NUMERICAL RESULTS

We numerically validated the theoretical results through simulation, and evaluated the performance of the simplified version of ANG D. The numerical results obtained using the theory were comparable with those of the simulation with respect to not only the learning curves but also the learning failures. We also found that the performance of ANG D is roughly comparable with that of NGD when  $\frac{\eta}{\rho}$  is small. Detailed conditions of these numerical results are given in Appendix D.

First, we validated the theoretical motion equations by using a simulation with  $N$  set to 500. The learning curves (time evolution of the generalization error) of the theory are shown in Fig. 2(c), while those of the simulation are shown in Fig. 3(c). There were no significant differences between the theoretical and the simulation results. With respect to the *two  $\xi$  rule*, simulation results showed that those adopting this rule were generally comparable with those not adopting this rule, although they were slightly slower when  $\rho$  is large. (We show only the simulation results for adoption of this rule.)

We also evaluated the learning failure of the simplified version of ANG D. As this version of ANG D defined by Eq. (12) assumes  $\rho \ll 1$ , a large  $\rho$  could cause problems. The simulation results showed divergence of the network parameters

and the system failed to learn the teacher outputs with large  $\rho$  (solid line in Fig. 4(a)). Figure 4(b) shows the borderline between this learning failure and the success areas with respect to the  $\eta$  and  $\rho$  conditions. Roughly, the failure area corresponded to  $\rho \geq 0.05$ . Our numerical solution of the theory successfully reproduced these learning failures in the simulation, which are shown as the dotted lines in Fig. 4(a) and Fig. 4(b). (We considered a learning failure to have occurred with the theory when the correlation matrix of the weight vectors,  $\mathbf{C}$ , violated its positive definiteness.)

Next, we compared the learning curves between SGD, NGD, and ANG D under various teacher weight vector correlations. (The angle between teacher weight vectors is denoted as  $\kappa$ .) Figure 2 shows the learning curves for (a) SGD, (b) NGD, and (c) ANG D. We can see that ANG D had almost the same performance as NGD and does not have any severe plateaus. Moreover, ANG D was not greatly affected by the teacher correlations, although SGD was.

Finally, we reveal the key condition affecting the learning plateau in ANG D. NGD is known to have a plateau when the learning rate  $\eta$  is too large [6]. We found that a plateau occurs in ANG D not only in the large  $\eta$  case, but also in the small  $\rho$  case. Figure 5(a) shows the time cost of learning under a wide range of  $\eta$  and  $\rho$ . This contour graph suggests that a plateau occurs when  $\frac{\eta}{\rho}$  is large. Our simulation study also supported this finding (Fig. 5(b)). This phenomenon may be interpreted to mean that  $\hat{\mathbf{G}}^{-1}$  cannot follow a change in the true  $\mathbf{G}^{-1}$  if  $\eta$  is relatively large compared to  $\rho$ .

#### V. CONCLUSION

We have developed a new order parameter expression for a simplified version of adaptive natural gradient learning in which the learning dynamics can be expressed using only a few order parameters. We numerically validated this theory through simulation and confirmed that this theory successfully reproduces not only the learning curve, but also the learning failure. We found that the ANG D performance is generally comparable with that of NGD. We also found that we can avoid the plateau in ANG D by making the update rate of the network parameter  $\eta$  low enough compared to the update rate of the inverse of the Fisher information matrix  $\rho$ .

#### Acknowledgments

This work was partially supported by Grant-in-Aid for Scientific Research on Priority Areas No. 14084212 and Grant-in-Aid for Scientific Research (C) No. 14580438.

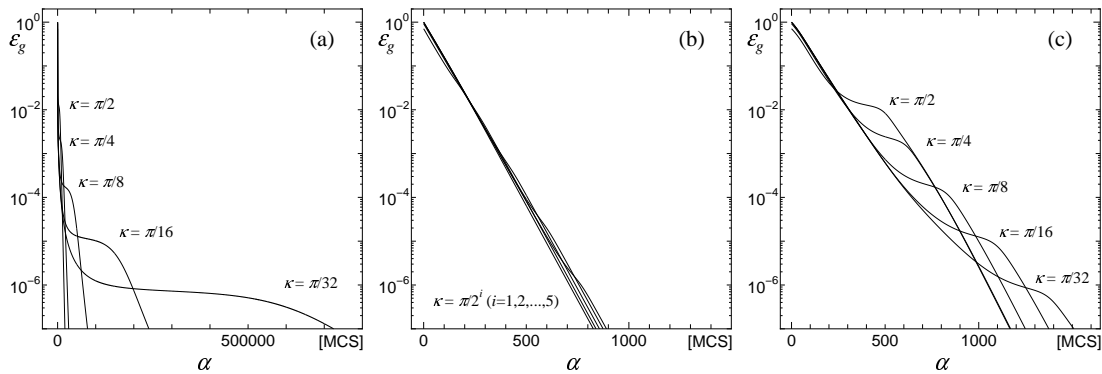


FIG. 2: Numerical results of the theory. Time evolution of the generalization error at  $\eta = 0.01, \rho = 0.01$ , and  $N = 500$ . (a) SGD, (b) NGD, and (c) ANGD. NGD and ANGD are not greatly affected by the angle of the teacher weight vectors ( $\kappa$ ), whereas SGD is.

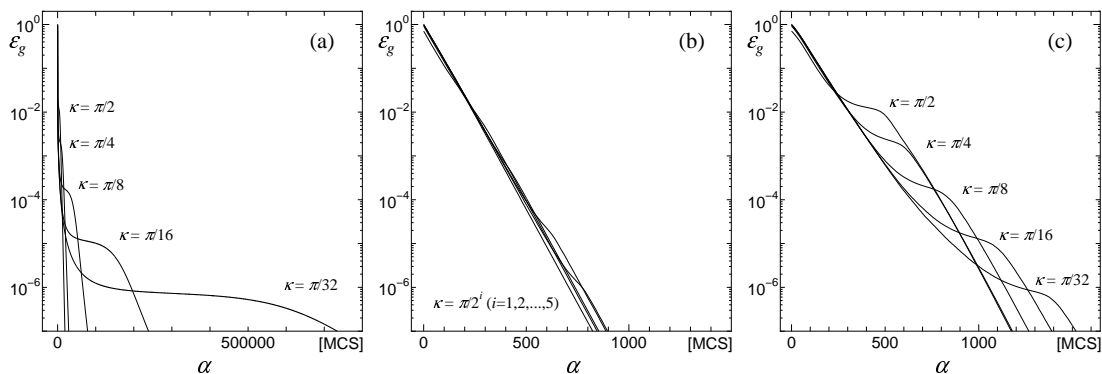


FIG. 3: Simulation results. Time evolution of the generalization error at  $\eta = 0.01, \rho = 0.01$ , and  $N = 500$ . (a) SGD, (b) NGD, and (c) ANGD. The results obtained using the theory in Fig. 2 are comparable to the simulation results.

## APPENDIX A: CONVERGENCE OF MATRICES

In this paper, we sometimes refer to matrix or vector ‘convergence’ in the sense of each element, although this word is usually used in the sense that the norm of the difference between the series of concern and a given matrix or vector converges to zero. As we are dealing with the large limit of the input dimension  $N$ , the average of  $N$   $\mathbb{R}^{N \times N}$ -matrices or  $N$   $\mathbb{R}^N$ -vectors often converges in the sense of each element, but does not in the sense of the norm. We can see one example of this phenomenon in  $\mathbf{J}_i$  dynamics in SGD; although each element of its fluctuation is small, the norm of this fluctuation is not zero but  $\Delta\phi_{ii}$  (see Eq. (21) and Fig. 1). In the following, we discuss the convergence of matrix  $\mathbf{F}$  used in Eq. (44).

Let us estimate the order of  $\mathbf{F}$ , specifically, the following matrix:

$$\frac{1}{N} \sum_{\mu=1}^N [\nabla f^{(\mu)}][\nabla f^{(\mu)}]^T, \quad (\text{A1})$$

where  $\nabla f^{(\mu)} = [g'(\mathbf{J}_i^T \boldsymbol{\xi}^{(\mu)}) \boldsymbol{\xi}^{(\mu)}]_{i=1,\dots,K} \in \mathbb{R}^{NK}$  is a random vector dependent on the random input  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Here, the superscript of  $\bullet^{(\mu)}$  denotes not the time but simple identification. The diagonal element of each block of this matrix, e.g.,

the  $(k, k)$  element of the  $(i, j)$  block, is

$$\frac{1}{N} \sum_{\mu=1}^N g'(\mathbf{J}_i^T \boldsymbol{\xi}^{(\mu)}) g'(\mathbf{J}_j^T \boldsymbol{\xi}^{(\mu)}) \xi_k^{(\mu)2}. \quad (\text{A2})$$

If we drop  $g'(\mathbf{J}_i^T \boldsymbol{\xi}^{(\mu)}) g'(\mathbf{J}_j^T \boldsymbol{\xi}^{(\mu)})$ , because it is  $O(1)$ , we notice that the probability distribution of this element is given by an  $N$ -freedom chi-square distribution; i.e., its moment-generating function is defined as

$$\varphi(t) = \left\{ 1 - 2 \left\{ \frac{1}{N} t \right\} \right\}^{-N/2}. \quad (\text{A3})$$

Consequently, we get the variance of this element as

$$\left. \frac{\partial^2}{\partial t^2} \ln \varphi(t) \right|_{t=0} = \frac{2}{N}. \quad (\text{A4})$$

The non-diagonal element of each block, e.g., the  $(k, l)$  element of the  $(i, j)$  block, is

$$\frac{1}{N} \sum_{\mu=1}^N g'(\mathbf{J}_i^T \boldsymbol{\xi}^{(\mu)}) g'(\mathbf{J}_j^T \boldsymbol{\xi}^{(\mu)}) \xi_k^{(\mu)} \xi_l^{(\mu)}. \quad (\text{A5})$$

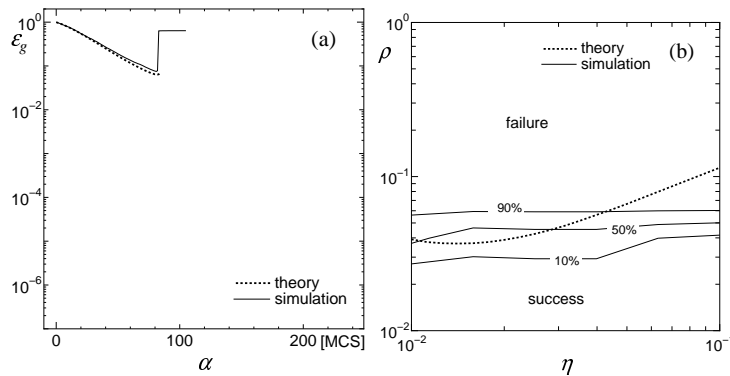


FIG. 4: Learning failure of the simplified version of ANG D with  $N = 100$ . (a) Time evolution of the generalization error under the conditions of  $\eta = 0.02$  and  $\rho = 0.1$ . (b) Contour graph under various values of  $\eta$  and  $\rho$ . The probability of learning failure is shown with respect to the simulation results because the simulation was a stochastic process. Numerical results of the theory well predicted the learning failure in the simulation.

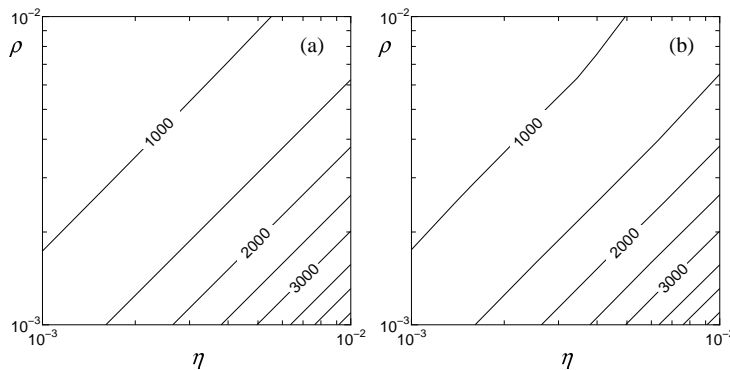


FIG. 5: Contour graphs of the time cost needed for  $\epsilon_g$  to reach  $\epsilon_g < 10^{-7}$  for various values of  $\eta$  and  $\rho$  in ANG D. Here,  $\kappa$  was set to  $\pi/8$ . (a) Results obtained using the theory, and (b) simulation results. Time  $\alpha$  was normalized as  $100\eta\alpha$ . The plateau length was strongly dependent on  $\frac{\eta}{\rho}$ .

Similarly, we notice that this distribution is given by an average of  $N$  modified Bessel functions of the second kind. Then, we get the moment generating function,

$$\varphi(t) = \left\{ 1 - \left\{ \frac{1}{N} t \right\}^2 \right\}^{-N/2}, \quad (\text{A6})$$

and the variance as  $\frac{1}{N}$ .

Therefore, each element of the matrix converges as  $O(\frac{1}{\sqrt{N}})$ , but the Frobenius norm diverges as  $O(\sqrt{N})$  because this matrix has  $NK \times NK$  elements. In other words,  $\sum F \rightarrow G$  holds with respect to each element, but does not converge with respect to the Frobenius norm.

## APPENDIX B: ORDERS OF $\omega, \Omega, \Omega', \Omega''$ AND $\Omega'''$

We prove that, with respect to  $N$ , the orders of each element of the matrices  $\omega, \Omega, \Omega', \Omega''$ , and  $\Omega'''$  are at most  $O(1), O(1), O(\frac{1}{\sqrt{N}}), O(\frac{1}{\sqrt{N}})$ , and  $O(\frac{1}{N})$ . First, we evaluate the

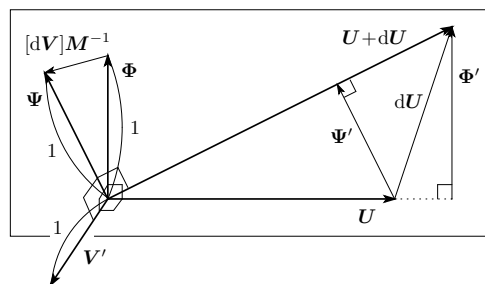


FIG. 6: Intuitive schema of the null space of  $U$  and  $U+dU$ .

order of the  $U^T dV$ ,  $V^T dU$ , and  $V^T dV$  terms, which are used in Eqs. (70)-(73). We especially pay attention to the fact that some of the inner products between infinitesimal changes of  $N$ -dimensional vectors are not  $O(d^2)$  but  $O(d)$ ; e.g.,  $[dU]^T dU = d\phi \equiv \begin{bmatrix} d\phi & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{(K+M) \times (K+M)}$ . We then consider Eqs. (69)-(73).

First of all, we explicitly determine  $V$ , the orthonormal

bases of the null space of the weight vectors, using  $U$  and  $dU$ . We find that some of the orthonormal bases of the  $V$  subspace can be expressed using  $U$  and  $dU$  as

$$\Phi \equiv \Phi' [\Phi'^T \Phi']^{-\frac{1}{2}} \in \mathbb{R}^{N \times (K+M)}, \quad (\text{B1})$$

where

$$\Phi' \equiv [I - UC^{-1}U^T] dU \in \mathbb{R}^{N \times (K+M)}. \quad (\text{B2})$$

Here,  $\Phi'$  is the orthogonal component of the  $dU$  to  $U$  subspace, while  $\Phi$  is the normalized  $\Phi'$ ; i.e.,  $\Phi^T \Phi = I$  (see Fig. 6). (We assume the rank of  $dU$  is  $K+M$  for simplicity, although it is actually  $K$ .) Note that  $UC^{-1}U^T$  is a projection matrix to the  $U$  subspace. In a similar manner, we also find that some of the orthonormal bases of the  $V+dV$  subspace can be expressed using  $U$  and  $dU$  as

$$\Psi \equiv \Psi' [\Psi'^T \Psi']^{-\frac{1}{2}} \in \mathbb{R}^{N \times (K+M)}, \quad (\text{B3})$$

where

$$\Psi' \equiv -[I - [U + dU][C + dC]^{-1}[U + dU]^T] U \in \mathbb{R}^{N \times (K+M)}. \quad (\text{B4})$$

Here,  $\Psi'$  is the orthogonal component of  $U$  to the  $U+dU$  subspace, while  $\Psi$  is the normalized  $\Psi'$ ; i.e.,  $\Psi^T \Psi = I$ . Thus, we can think that some of the orthonormal bases of the  $V$  subspace correspond to  $\Phi$ , and  $\Phi$  moves to  $\Psi$  as  $V$  moves to  $V+dV$ . However, the column vectors of  $\Phi$  and  $\Psi$  do not necessarily coincide with some of the column vectors of  $V$  and  $V+dV$ , respectively; i.e., some rotation or mirror image conversion might be required. Hence, we introduce an appropriate orthonormal matrix  $M \in \mathbb{R}^{(N-K-M) \times (N-K-M)}$ , and explicitly express  $V$ ,  $V+dV$ , and  $dV$  as

$$V = [\Phi, V']M, \quad (\text{B5})$$

$$V+dV = [\Psi, V']M, \quad (\text{B6})$$

$$dV = [\Psi - \Phi, \mathbf{0}]M, \quad (\text{B7})$$

where  $V'$  is one of the set of orthonormal bases of the null space of  $U$  and  $U+dU$ . Note that each of

$$\begin{aligned} & [UC^{-\frac{1}{2}}, V], \\ & [UC^{-\frac{1}{2}}, \Phi, V'], \\ & [U+dU][C+dC]^{-\frac{1}{2}}, V+dV], \\ & [U+dU][C+dC]^{-\frac{1}{2}}, \Psi, V'], \end{aligned} \quad (\text{B8})$$

consists of the orthonormal bases of the whole space. If we decompose  $M$  as  $\begin{bmatrix} M_0 \\ M_1 \end{bmatrix}$ , we can rewrite Eq. (B5)-(B7) as

$$V = [\Phi, V']M = \Phi M_0 + V' M_1, \quad (\text{B9})$$

$$V+dV = [\Psi, V']M = \Psi M_0 + V' M_1, \quad (\text{B10})$$

$$dV = [\Psi - \Phi, \mathbf{0}]M = [\Psi - \Phi]M_0. \quad (\text{B11})$$

The  $V'$  and  $M_1$  contain some arbitrariness, although  $V' M_1$  is well-defined. However, we can avoid using them as shown below.

We can then calculate the  $V^T dU$ ,  $[dV]^T U$ , and  $[dV]^T V$  terms. We find that the norms of the column vectors of  $V^T dU$ ,

$$\begin{aligned} & [dU]^T V V^T dU \\ & = [dU]^T [\Phi, V'] M M^T [\Phi, V']^T dU \\ & = d\phi', \end{aligned} \quad (\text{B12})$$

are  $O(1)$  with respect to  $N$ . Note that  $M M^T = I$  by its definition. As the orthonormal matrix  $M$  does not change a matrix norm, each element of the  $N$ -dimensional vectors  $V^T dU$  is  $O(\frac{1}{\sqrt{N}})$ . In a similar manner, we can calculate the norm of  $[dV]^T U$ ,

$$\begin{aligned} & U^T [dV][dV]^T U \\ & = U^T [\Psi - \Phi, \mathbf{0}] M M^T [\Psi - \Phi, \mathbf{0}]^T U \\ & = d\phi'. \end{aligned} \quad (\text{B13})$$

This also means that each element of  $[dV]^T U$  is  $O(\frac{1}{\sqrt{N}})$ . These results are used after Eq. (73). Equation (B13) is also used to derive Eq. (82). We also find that each element of  $[dV]^T V$  is  $O(\frac{1}{\sqrt{N^2}})$ , because

$$\begin{aligned} & V^T [dV] \\ & = M^T [\Phi, V']^T [\Psi - \Phi, \mathbf{0}] M \\ & = -M^T \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} M \\ & = -M_0^T A M_0, \end{aligned} \quad (\text{B14})$$

where each element of  $A \equiv [d\phi']^{-\frac{1}{2}} [[d\phi']C^{-1}[dC-dO] + [dO]C^{-1}dO^T] [d\phi']^{-\frac{1}{2}}$  is  $O(1)$ , and these elements are scattered to  $N \times N$  elements by  $M$ .

Next, we consider the orders of Eqs. (69)-(73):

$$d\omega_{ij} = \rho[\omega - \omega\lambda\omega]_{ij} d\alpha, \quad (\text{B15})$$

$$\begin{aligned} d\Omega_{ij} = & \rho C^{-\frac{1}{2}} [\tilde{\Omega} - \tilde{\Omega}\tilde{\Lambda}\tilde{\Omega} - \tilde{\omega} + \tilde{\omega}\tilde{\lambda}\tilde{\omega} - \tilde{\Omega}'\tilde{\lambda}\tilde{\Omega}'']_{ij} C^{-\frac{1}{2}} d\alpha \\ & - C^{-1}[dC-dO]\Omega_{ij} - \Omega_{ij}[dC-dO]^T C^{-1} + C^{-1}[dU]^T V\Omega_{ij}'' + \Omega_{ij}'' V^T [dU]C^{-1}, \end{aligned} \quad (\text{B16})$$

$$\begin{aligned} d\Omega'_{ij} &= \rho C^{-\frac{1}{2}} \left[ \tilde{\Omega}' - \tilde{\Omega} \tilde{\Lambda} \tilde{\Omega}' - \tilde{\Omega}' \tilde{\lambda} \tilde{\omega} - \tilde{\Omega}' \tilde{\lambda} \Omega''' \right]_{ij} d\alpha \\ &\quad - C^{-1} [dC - dO] \Omega'_{ij} + \Omega'_{ij} V^T [dV] + C^{-1} [dU]^T V \Omega''_{ij} + \Omega_{ij} U^T [dV], \end{aligned} \quad (B17)$$

$$\begin{aligned} d\Omega''_{ij} &= \rho \left[ \tilde{\Omega}'' - \tilde{\Omega}'' \tilde{\Lambda} \tilde{\Omega}'' - \tilde{\omega} \tilde{\lambda} \tilde{\Omega}'' - \Omega''' \tilde{\lambda} \tilde{\Omega}'' \right]_{ij} C^{-\frac{1}{2}} d\alpha \\ &\quad + [dV]^T V \Omega''_{ij} - \Omega''_{ij} [dC - dO]^T C^{-1} + [dV]^T U \Omega_{ij} + \Omega''_{ij} V^T [dU] C^{-1}, \end{aligned} \quad (B18)$$

$$\begin{aligned} d\Omega'''_{ij} &= \rho \left[ \Omega''' - \tilde{\Omega}'' \tilde{\Lambda} \tilde{\Omega}' - \tilde{\omega} \tilde{\lambda} \Omega''' - \Omega''' \tilde{\lambda} \tilde{\omega} - \Omega''' \tilde{\lambda} \Omega''' \right]_{ij} d\alpha \\ &\quad + [dV]^T V \Omega'''_{ij} + \Omega'''_{ij} V^T [dV] + [dV]^T U \Omega'_{ij} + \Omega'''_{ij} U^T [dV]. \end{aligned} \quad (B19)$$

Each element of the matrices  $C$ ,  $dC$ , and  $dO$  are considered to be  $O(1)$  with respect to  $N$ . The  $\lambda$  and  $\mathbf{A}$  are also considered to be  $O(1)$ . Matrix normalization denoted with tilde (e.g.,  $\tilde{\Omega}$ ) is considered to not change the order. Matrix size extension (e.g.,  $\tilde{\omega}$ ) is also considered to not change the order. As the initial value of  $\mathbf{H}$  is defined as a unit matrix, we can let  $\omega = \mathbf{I}$ ,  $\Omega = \mathbf{0}$ ,  $\Omega' = \mathbf{0}$ ,  $\Omega'' = \mathbf{0}$ , and  $\Omega''' = \mathbf{0}$  as initial values. Then, we notice the following. 1)  $\omega$  is  $O(1)$  from the initial state. 2) Then,  $\Omega$  soon becomes  $O(1)$ , because  $d\Omega_{ij}$  has a  $O(1)$  term:  $-\tilde{\omega} + \tilde{\omega} \tilde{\lambda} \tilde{\omega}$ . 3) Then,  $\Omega'$  and  $\Omega''$  soon become  $O(\frac{1}{\sqrt{N}})$ , because  $d\Omega'_{ij}$  and  $d\Omega''_{ij}$  have a  $O(\frac{1}{\sqrt{N}})$  term:  $\Omega_{ij} U^T [dV]$  and  $[dV]^T U \Omega_{ij}$ , respectively. 4) Then,  $\Omega'''$  soon becomes  $O(\frac{1}{\sqrt{N}})$ , because  $d\Omega'''_{ij}$  has a  $O(\frac{1}{\sqrt{N}})$  term:  $[dV]^T U \Omega'_{ij} + \Omega'''_{ij} U^T [dV]$ . 5) There is no contradiction if we assume these orders are preserved.

### APPENDIX C: EFFECT OF $E$

We find that the fluctuation term  $E$  defined by Eq. (50) as

$$E_{ij} \equiv [U, V] \begin{bmatrix} \mathbf{0} & \Omega'_{ij} \\ \Omega''_{ij} & \Omega'''_{ij} \end{bmatrix} [U, V]^T, \quad (C1)$$

does not affect the usual order parameter dynamics under the assumption that each element of  $\Omega'_{ij}$ ,  $\Omega''_{ij}$  and  $\Omega'''_{ij}$  are  $O(\frac{1}{\sqrt{N}})$ ,  $O(\frac{1}{\sqrt{N}})$ , and  $O(\frac{1}{\sqrt{N}})$ , respectively (see Appendix B).

Here,  $E$  is negligible with respect to the dynamics of  $\mathbf{R}$  and  $\mathbf{S}$  in Eq. (47). This is because all the terms including  $E$  are 0 as shown

$$\begin{aligned} \langle \delta_k \xi^T E_{ik}^T \mathbf{B}_j \rangle_{\xi} &= \langle \delta_k \xi^T V \Omega''_{ik}{}^T U^T \mathbf{B}_j \rangle_{\xi} \\ &= \langle \check{\epsilon}(z) g'(x_k) \rangle_z \langle v^T \rangle_v \Omega''_{ik}{}^T U^T \mathbf{B}_j = 0, \\ \langle \delta_k \xi^T E_{ik}^T \mathbf{J}_j \rangle_{\xi} &= \langle \delta_k \xi^T V \Omega''_{ik}{}^T U^T \mathbf{J}_j \rangle_{\xi} \\ &= \langle \check{\epsilon}(z) g'(x_k) \rangle_z \langle v^T \rangle_v \Omega''_{ik}{}^T U^T \mathbf{J}_j = 0, \end{aligned} \quad (C2)$$

where  $v \equiv V^T \xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is independent of  $z \equiv U^T \xi = [x_1, \dots, x_K, y_1, \dots, y_M]^T$ , and  $\langle v \rangle_v = \mathbf{0}$ .

$E$  is also negligible with respect to the dynamics of  $\phi$ . All the terms including  $E$  can be expressed as

$$\frac{1}{N} \left\langle \delta_k \delta_l \xi^T [U, V] \begin{bmatrix} \mathbf{A} & \mathbf{A}' \\ \mathbf{A}'' & \mathbf{A}''' \end{bmatrix} [U, V]^T \xi \right\rangle_{\xi}, \quad (C3)$$

where

$$\begin{bmatrix} \mathbf{A} & \mathbf{A}' \\ \mathbf{A}'' & \mathbf{A}''' \end{bmatrix} = \mathbf{H}_{ik} \mathbf{H}_{jl} - [\omega_{ik} + U \Omega_{ik} U^T][\omega_{jl} + U \Omega_{jl} U^T] \quad (C4)$$

Then, we notice that  $\langle \delta_k \delta_l \xi^T U \mathbf{A} U \xi \rangle_{\xi}$  is  $O(1)$ , while  $\langle \delta_k \delta_l \xi^T U \mathbf{A}' V \xi \rangle_{\xi}$  and  $\langle \delta_k \delta_l \xi^T V \mathbf{A}'' U \xi \rangle_{\xi}$  are 0. Moreover,  $\langle \delta_k \delta_l \xi^T V \mathbf{A}''' V \xi \rangle_{\xi}$  is at most  $O(1)$  because each element of  $\mathbf{A}'''$ ,

$$\mathbf{A}''' = \omega_{ik} \Omega''_{jl} + \omega_{jl} \Omega''_{ik} + \Omega''_{ik} C \Omega'_{jl} + \Omega''_{ik} C \Omega'''_{ik}, \quad (C5)$$

is  $O(\frac{1}{\sqrt{N}})$ . As all the terms in Eq. (C3) are scaled by  $\frac{1}{\sqrt{N}}$ , we can ignore them. Therefore, we can completely ignore the effect of  $E$  in the usual order parameter dynamics.

### APPENDIX D: DETAILED CONDITIONS FOR NUMERICAL RESULTS

For numerical results, we considered a realizable case, in which the numbers of the hidden units for both the teacher and student networks were set to two ( $K = M = 2$ ). With respect to the order parameter dynamics, the initial conditions of the usual order parameters were set as follows. The square lengths of all teacher weight vectors  $T_{ii}$  were set to 1, while the angle between the teacher weight vectors,  $\kappa \equiv \arccos \frac{T_{1,2}}{\sqrt{T_{1,1} T_{2,2}}}$ , was set to a moderately correlated value,  $\pi/8$ , unless otherwise stated. The initial conditions with respect to the student weight vectors were determined according to the corresponding expected values of random choice  $\mathbf{J} \sim \mathcal{N}(\mathbf{0}, \frac{1}{N} \mathbf{I})$ ; i.e.,  $Q_{ii} = 1$ ,  $Q_{ij} = 0$ , and  $R_{ii} = 0$ . Only  $R_{ij} (i \neq j)$ , the inner products between the student weight vectors and non-corresponding teacher vectors, were set to a small negative value  $R_{ii-r_{ij}}$  to break the permutation symmetry a little, where we adopted a 1 standard deviation (S.D.) rule; i.e.,

$$r_{ij} \equiv \sqrt{\text{Var}_{\mathbf{J}}(R_{ii} - R_{ij})} = \sqrt{\frac{T_{ii} + T_{jj} - 2T_{ij}}{N}}, \quad (D1)$$

where  $\text{Var}_{\mathbf{J}}(\bullet)$  denotes the variance of  $\bullet$  with respect to  $\mathbf{J}$ . Then, we solved the order parameter equations using the Runge-Kutta method with time interval  $\Delta\alpha = 0.1$ . With respect to our simulation, the initial states of  $\mathbf{B}$  and  $\mathbf{J}$  were determined to satisfy the conditions of the order parameters above.

- 
- [1] D. Saad and S. A. Solla: Phys. Rev. Lett. **74** 4337 (1995).  
D. Saad and S. A. Solla: Phys. Rev. E **52** 4225 (1995).
- [2] P. Riegler and M. Biehl: J. Phys. A **28** L507 (1995).
- [3] M. Inoue, H. Park and M. Okada: J. Phys. Soc. Jpn. **72** 805 (2003).
- [4] S. Amari: Neural Comput. **10** 251 (1998).
- [5] H. H. Yang and S. Amari: Neural Comput. **10** 2137 (1998).
- [6] M. Rattray, D. Saad and S. Amari: Phys. Rev. Lett. **81** 5461 (1998).  
M. Rattray and D. Saad: Phys. Rev. E **59** 4523 (1999).
- [7] Online Learning in Neural Networks, edited by D. Saad, Cambridge University Press, London (1998).
- [8] S. Amari, H. Park and K. Fukumizu: Neural Comput. **12** 1399 (2000).
- [9] H. Park, S. Amari and K. Fukumizu: Neural Networks **13** 755 (2000).
- [10] J. Inoue, H. Nishimori and Y. Kabashima: J. Phys. A **30** 3795 (1997).
- [11] D. Saad and M. Rattray: Physical Review Letters **79** 2578 (1997).  
M. Rattray and D. Saad: Physical Review E **58**, 6379 (1998).