

Natural Gradient Descent for Training Stochastic Complex-Valued Neural Networks

Tohru Nitta

National Institute of Advanced Industrial Science and Technology (AIST),
AIST Tsukuba Central 2, 1-1-1 Umezono, Tsukuba, Ibaraki 305-8568, Japan
Email: tohru-nitta@aist.go.jp

Abstract—In this paper, the natural gradient descent method for the multilayer stochastic complex-valued neural networks is considered, and the natural gradient is given for a single stochastic complex-valued neuron as an example. Since the space of the learnable parameters of stochastic complex-valued neural networks is not the Euclidean space but a curved manifold, the complex-valued natural gradient method is expected to exhibit excellent learning performance.

Keywords—Neural network; Complex number; Learning; Singular point

I. INTRODUCTION

Complex-valued neural networks whose parameters (weights and threshold values) are all complex numbers, are useful in fields dealing with complex numbers or two-dimensional vectors such as telecommunications, speech recognition and image processing with Fourier transformation. Indeed, we can find some applications of complex-valued neural networks to various fields in the literature [6], [9].

The multilayer complex-valued neural network is usually trained using the gradient descent learning method [5], [10], [11], [12], as in the case of the multilayer real-valued neural network. The space of the learnable parameters of stochastic complex-valued neural networks is, however, not the Euclidean space but a curved manifold. For stochastic complex-valued neural networks, the ordinary gradient does not give the steepest direction of a target function, and the steepest direction is given by the natural gradient [2], [3]. It has been shown in [4] that the natural gradient method could avoid singular points of the real-valued parameter space which is a cause of standstill in learning, and the natural gradient method could improve the learning performance of the real-valued neural networks as a result. Similarly, there exist many singular points in the complex-valued neural networks [7]. Thus, the natural gradient method would be useful for the complex-valued neural networks, too. In this paper, we extend the natural gradient descent method for the multilayer stochastic real-valued neural networks to the complex domain, and give the natural gradient for a single stochastic complex-valued neuron as an example.

Section II describes the complex-valued neural network. Section III is devoted to the explanation of the natural gradient method, and Section IV presents the natural gradient method in complex-valued neural networks, which is followed by our conclusion in Section V.

II. COMPLEX-VALUED NEURAL NETWORK MODEL

This section describes the complex-valued neural network model used in this paper. First, we will consider the following complex-valued neuron. The input signals, weights, thresholds and output signals are all complex numbers. The net input U_n to a complex-valued neuron n is defined as: $U_n = \sum_m W_{nm} X_m + V_n$, where W_{nm} is the complex-valued weight connecting the complex-valued neurons n and m , X_m is the complex-valued input signal from the complex-valued neuron m , and V_n is the complex-valued threshold value of the complex-valued neuron n . To obtain the complex-valued output signal, convert the net input U_n into its real and imaginary parts as follows: $U_n = x + iy = z$, where i denotes $\sqrt{-1}$. The complex-valued output signal is defined to be

$$f_C(z) = \varphi(x) + i\varphi(y), \quad (1)$$

where $\varphi: \mathbf{R} \rightarrow \mathbf{R}$, (\mathbf{R} denotes the set of real numbers). Eq. (1) is often called a *split-type* complex-valued activation function. Note that the activation function f_C is not a regular complex-valued function because the Cauchy-Riemann equations do not hold.

The complex-valued neural network used in this paper consists of such complex-valued neurons described above.

Note that various types of activation functions other than Eq. (1) can be considered naturally (for examples, the non-split-type (fully) one [10]).

III. NATURAL GRADIENT METHOD

This section briefly describes the natural gradient proposed in [2], [3]. Let $S = \{\mathbf{w} \in \mathbf{R}^N\}$ be a Riemannian space with the Riemannian metric tensor $\mathbf{G}(\mathbf{w}) = (g_{ij}(\mathbf{w}))$ on which a function $L(\mathbf{w})$ is defined. If

$$g_{ij}(\mathbf{w}) = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases}, \quad (2)$$

that is, $\mathbf{G}(\mathbf{w})$ is the unit matrix, then S is an Euclidean space. Amari proved the following theorem [3].

Theorem 1: *The steepest descent direction of $L(\mathbf{w})$ in a Riemannian space is given by*

$$-\tilde{\nabla}L(\mathbf{w}) = -\mathbf{G}^{-1}(\mathbf{w})\nabla L(\mathbf{w}) \quad (3)$$

where $\mathbf{G}^{-1}(\mathbf{w}) = (g^{ij}(\mathbf{w}))$ is the inverse of the metric $\mathbf{G}(\mathbf{w}) = (g_{ij}(\mathbf{w}))$ and ∇L is the conventional gradient,

$$\nabla L(\mathbf{w}) = \left(\frac{\partial}{\partial w_1} L(\mathbf{w}), \dots, \frac{\partial}{\partial w_N} L(\mathbf{w}) \right)^T, \quad (4)$$

where the superscript T denotes the transposition. \square

$\tilde{\nabla}L(\mathbf{w}) = \mathbf{G}^{-1}\nabla L(\mathbf{w})$ is called the *natural gradient* of L in the Riemannian space. The natural gradient descent algorithm is given by

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \varepsilon_t \tilde{\nabla}L(\mathbf{w}_t), \quad (5)$$

where ε_t is the learning rate.

Amari derived the natural gradients explicitly in the case of the space of real-valued perceptrons for neural learning, the space of matrices for blind source separation, and the space of linear dynamical systems for blind multichannel source deconvolution [3].

IV. NATURAL GRADIENT IN COMPLEX-VALUED NEURAL NETWORKS

In this section, the natural gradient is applied to the complex-valued neural networks and the natural gradient descent algorithm is explicitly derived for a single complex-valued neuron.

A. Natural Gradient Learning in Complex-Valued Neural Networks

Let us consider a stochastic complex-valued multilayer feedforward neural network with N input neurons, one output neuron, and a learnable complex-valued vector parameter $\mathbf{w} = (w_1, \dots, w_N)^T \in \mathbf{C}^N$ which consists of all the weights and thresholds (\mathbf{C} denotes the set of complex numbers). Assume that the complex-valued input signal $\mathbf{z} = (z_1, \dots, z_N)^T \in \mathbf{C}^N$ is subject to an unknown probability distribution $q(\mathbf{z})$, and the complex-valued output signal $y \in \mathbf{C}$ is given by

$$y = g_C(\mathbf{z}, \mathbf{w}) + n, \quad (6)$$

where g_C is a complex function, and $n = n_R + in_I$ is a complex-valued random variable subject to a complex normal distribution (or bivariate normal distribution) $N(\mu, \Sigma)$. The model specifies the probability density of the input-output pair as

$$p(\mathbf{z}, y; \mathbf{w}) = q(\mathbf{z}) \cdot p(y|\mathbf{z}; \mathbf{w}). \quad (7)$$

Define a loss function $l(\mathbf{z}, y; \mathbf{w})$ when input signal \mathbf{z} is processed by the stochastic complex-valued neural network having parameter \mathbf{w} as:

$$\begin{aligned} l(\mathbf{z}, y; \mathbf{w}) &\stackrel{\text{def}}{=} -\log p(\mathbf{z}, y; \mathbf{w}) \\ &= -\log q(\mathbf{z}) - \log p(y|\mathbf{z}; \mathbf{w}). \end{aligned} \quad (8)$$

Given the training set $\{(z_t, y_t), t = 1, \dots, T\}$, minimizing the loss function (Eq. (8)) is equivalent to maximizing the probability that the stochastic complex-valued neural network outputs the training output signal y_t .

The space of all the probability distributions which the above stochastic complex-valued neural network realizes, can be regarded as a $2N$ -dimensional Riemannian space because the complex-valued parameter consists of the two real-valued parameters: the real-part and the imaginary part. Thus, the information geometry [1] can be applied to the complex-valued case, too.

The natural gradient descent algorithm for the complex-valued neural network is given by

$$\mathbf{v}_{t+1} = \mathbf{v}_t - \varepsilon_t \tilde{\nabla}l(\mathbf{z}_t, y_t; \mathbf{v}_t), \quad (9)$$

where $\{(z_t, y_t) \in \mathbf{C}^N \times \mathbf{C}, t = 1, 2, \dots\}$ is the sequence of the complex-valued training signals, and

$$\begin{aligned} \mathbf{v} &= (v_1, \dots, v_{2N})^T \\ &= (\text{Re}[w_1], \dots, \text{Re}[w_N], \text{Im}[w_1], \dots, \text{Im}[w_N])^T, \end{aligned} \quad (10)$$

$$\tilde{\nabla}l(\mathbf{z}, y, \mathbf{v}) \stackrel{\text{def}}{=} \mathbf{G}^{-1}(\mathbf{v}) \cdot \nabla l(\mathbf{z}, y, \mathbf{v}). \quad (11)$$

Eq. (11) is the natural gradient of $l(\mathbf{z}, y, \mathbf{v})$, and the usual gradient $\nabla l(\mathbf{z}, y, \mathbf{v})$ is given by

$$\begin{aligned} \nabla l(\mathbf{z}, y, \mathbf{v}) &\stackrel{\text{def}}{=} \left(\frac{\partial l(\mathbf{z}, y, \mathbf{v})}{\partial v_1}, \dots, \frac{\partial l(\mathbf{z}, y, \mathbf{v})}{\partial v_{2N}} \right)^T \\ &= \left(\frac{\partial l(\mathbf{z}, y, \mathbf{w})}{\partial \text{Re}[w_1]}, \dots, \frac{\partial l(\mathbf{z}, y, \mathbf{w})}{\partial \text{Re}[w_N]}, \right. \\ &\quad \left. \frac{\partial l(\mathbf{z}, y, \mathbf{w})}{\partial \text{Im}[w_1]}, \dots, \frac{\partial l(\mathbf{z}, y, \mathbf{w})}{\partial \text{Im}[w_N]} \right)^T. \end{aligned} \quad (12)$$

The Riemannian metric tensor $\mathbf{G}(\mathbf{v})$ is the Fisher information matrix [3], and is given by

$$\mathbf{G}(\mathbf{v}) = (g_{ij}(\mathbf{v})), \quad (13)$$

$$g_{ij}(\mathbf{v}) = E \left[\frac{\partial \log p(\mathbf{z}, y; \mathbf{v})}{\partial v_i} \cdot \frac{\partial \log p(\mathbf{z}, y; \mathbf{v})}{\partial v_j} \right]. \quad (14)$$

B. Natural Gradient Learning in a Single Complex-Valued Neuron

In this section, the natural gradient descent learning algorithm for a single complex-valued neuron is given.

Consider a stochastic complex-valued neuron with N -inputs, weights $w_k = u_k + iv_k \in \mathbf{C}$ ($1 \leq k \leq N$), and a threshold value $\gamma = c + id \in \mathbf{C}$. Then, for N input signals $z_k = x_k + iy_k \in \mathbf{C}$ ($1 \leq k \leq N$), the stochastic complex-valued neuron generates

$$\begin{aligned} y &= f_C \left(\sum_{k=1}^N w_k z_k + \gamma \right) + n \\ &= X + iY \end{aligned} \quad (15)$$

where $f_C : \mathbf{C} \rightarrow \mathbf{C}$ is a so-called *split-type* complex-valued activation function which is defined to be

$$f_C(a + ib) = \varphi(a) + i\varphi(b) \quad (16)$$

for any $a + ib \in \mathbf{C}$, and $\varphi : \mathbf{R} \rightarrow \mathbf{R}$ is suitably chosen, for example, the sigmoid function

$$\varphi(s) = \frac{1}{1 + e^{-s}} \quad (17)$$

was used in [11], and the scaled error function

$$\varphi(s) = \frac{2}{\sqrt{\pi}} \int_0^{\frac{s}{\sqrt{2}}} e^{-t^2} dt \quad (18)$$

was used in [13]. $n = n_R + in_I$ is a complex-valued random variable subject to the complex normal distribution (or bivariate normal distribution) $N(\mu, \Sigma)$ where

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad (19)$$

$$\Sigma = \begin{bmatrix} \sigma & 0 \\ 0 & \sigma \end{bmatrix}. \quad (20)$$

We assume that the input signal $\mathbf{z} = \mathbf{x} + i\mathbf{y}$ where $\mathbf{x} \stackrel{\text{def}}{=} (x_1, \dots, x_N)^T, \mathbf{y} \stackrel{\text{def}}{=} (y_1, \dots, y_N)^T$ is subject to the multivariate complex normal distribution (or $2N$ -dimensional normal distribution) $N(\mathbf{0}, \mathbf{I})$ where the variance covariance matrix \mathbf{I} is the unit matrix; denote its joint probability density function by $q(\mathbf{z})$. The loss function $l(\mathbf{z}, y; \theta)$ is defined as

$$l(\mathbf{z}, y; \theta) \stackrel{\text{def}}{=} -\log p(\mathbf{z}, y; \theta) \\ = -\log q(\mathbf{z}) - \log p(y|\mathbf{z}; \theta), \quad (21)$$

where $\theta \stackrel{\text{def}}{=} (\mathbf{u}^T, \mathbf{v}^T, c, d)^T, \mathbf{u} = (u_1, \dots, u_N)^T, \mathbf{v} = (v_1, \dots, v_N)^T$.

Given the sequence of the complex-valued training signals $\{(\mathbf{z}_t, y_t) \in \mathbf{C}^N \times \mathbf{C}, t = 1, 2, \dots\}$, the the natural gradient descent algorithm for the stochastic complex-valued neuron is given by

$$\theta_{t+1} = \theta_t - \varepsilon_t \cdot \mathbf{G}^{-1}(\theta_t) \cdot \nabla l(\mathbf{z}_t, y_t; \theta_t). \quad (22)$$

We shall calculate the Fisher information matrix $\mathbf{G}(\theta) = (g_{ij}(\theta))$. For any $1 \leq i, j \leq 2N + 2$,

$$g_{ij}(\theta) = E \left[\frac{\partial \log p(\mathbf{z}, y; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(\mathbf{z}, y; \theta)}{\partial \theta_j} \right] \\ = E \left[\frac{\partial \log p(y|\mathbf{z}; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(y|\mathbf{z}; \theta)}{\partial \theta_j} \right]. \quad (23)$$

(from Eq. (21))

Here, since n_R is independent of n_I ,

$$\log p(y|\mathbf{z}; \theta) = \log p((\varphi(S) + n_R) + i(\varphi(T) + n_I)|\mathbf{z}; \theta) \\ = \log p(\varphi(S) + n_R|\mathbf{z}; \theta) \\ + \log p(\varphi(T) + n_I|\mathbf{z}; \theta) \\ = \log p(X|\mathbf{z}; \theta) + \log p(Y|\mathbf{z}; \theta), \quad (24)$$

where $S = \text{Re} \left[\sum_{k=1}^N w_k z_k + \gamma \right], T = \text{Im} \left[\sum_{k=1}^N w_k z_k + \gamma \right]$. Thus, for any $1 \leq i, j \leq 2N + 2$,

$$g_{ij}(\theta) = E \left[\frac{\partial \log p(X|\mathbf{z}; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(X|\mathbf{z}; \theta)}{\partial \theta_j} \right] \\ + 2E \left[\frac{\partial \log p(X|\mathbf{z}; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(Y|\mathbf{z}; \theta)}{\partial \theta_j} \right] \\ + E \left[\frac{\partial \log p(Y|\mathbf{z}; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(Y|\mathbf{z}; \theta)}{\partial \theta_j} \right]. \quad (25)$$

By simple calculations, we obtain

$$g_{ij}(\theta) = \frac{1}{\sigma^2} \{ E[(\varphi'(S))^2 x_i x_j] \\ + E[(\varphi'(T))^2 y_i y_j] \} \\ (1 \leq i, j \leq N), \quad (26)$$

$$g_{ij}(\theta) = \frac{1}{\sigma^2} \{ E[(\varphi'(S))^2 (-y_{i-N}) x_j] \\ + E[(\varphi'(T))^2 x_{i-N} y_j] \} \\ = g_{ji}(\theta) \\ (N + 1 \leq i \leq 2N, 1 \leq j \leq N), \quad (27)$$

$$g_{ij}(\theta) = \frac{1}{\sigma^2} \{ E[(\varphi'(S))^2 y_{i-N} y_{j-N}] \\ + E[(\varphi'(T))^2 x_{i-N} x_{j-N}] \} \\ (N + 1 \leq i, j \leq 2N), \quad (28)$$

$$g_{2N+1,j}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(S))^2 x_j] \\ = g_{j,2N+1}(\theta) \quad (1 \leq j \leq N), \quad (29)$$

$$g_{2N+1,j}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(S))^2 (-y_{j-N})] \\ (N + 1 \leq j \leq 2N), \quad (30)$$

$$g_{2N+2,j}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(T))^2 y_j] \quad (1 \leq j \leq N), \quad (31)$$

$$g_{2N+2,j}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(T))^2 x_{j-N}] \\ (N + 1 \leq j \leq 2N), \quad (32)$$

$$g_{2N+1,2N+1}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(S))^2], \quad (33)$$

$$g_{2N+2,2N+1}(\theta) = 0 \\ = g_{2N+1,2N+2}(\theta), \quad (34)$$

$$g_{2N+2,2N+2}(\theta) = \frac{1}{\sigma^2} E[(\varphi'(T))^2]. \quad (35)$$

From Eqs. (26) – (35), we can rewrite $\mathbf{G}(\theta)$ as

$$\mathbf{G}(\theta) = \frac{1}{\sigma^2} \mathbf{A}(\theta), \quad (36)$$

where

$$\mathbf{A}(\theta) \stackrel{\text{def}}{=} \left(\begin{array}{cc|cc} \mathbf{A}_{11} & \mathbf{A}_{12} & \mathbf{a}_{13} & \mathbf{a}_{14} \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \mathbf{a}_{23} & \mathbf{a}_{24} \\ \hline \mathbf{a}_{31} & \mathbf{a}_{32} & & \\ \mathbf{a}_{41} & \mathbf{a}_{42} & & \mathbf{A}_{44} \end{array} \right), \quad (37)$$

$\mathbf{A}_{11} = \sigma^2 \cdot \text{Eq. (26)}, \mathbf{A}_{21} = \mathbf{A}_{12}^T = \sigma^2 \cdot \text{Eq. (27)}, \mathbf{A}_{22} = \sigma^2 \cdot \text{Eq. (28)}, \mathbf{a}_{13} = \mathbf{a}_{31}^T = \sigma^2 \cdot \text{Eq. (29)}, \mathbf{a}_{32} = \mathbf{a}_{23}^T = \sigma^2 \cdot \text{Eq. (30)}, \mathbf{a}_{14} = \mathbf{a}_{41}^T = \sigma^2 \cdot \text{Eq. (31)}, \mathbf{a}_{42} = \mathbf{a}_{24}^T = \sigma^2 \cdot \text{Eq. (32)},$ and

$$\mathbf{A}_{44} = \left(\begin{array}{cc} \sigma^2 \cdot \text{Eq. (33)} & \text{Eq. (34)} \\ \text{Eq. (34)} & \sigma^2 \cdot \text{Eq. (35)} \end{array} \right). \quad (38)$$

In what follows, each submatrix of $\mathbf{A}(\theta)$ (Eq. (37)) is calculated. Let $u = \|\mathbf{u}\| = \sqrt{u_1^2 + \dots + u_N^2}, v = \|\mathbf{v}\| = \sqrt{v_1^2 + \dots + v_N^2}, \mathbf{u}_1 = \mathbf{u}/u, \mathbf{v}_1 = \mathbf{v}/v$, and extend $\mathbf{u}_1, \mathbf{v}_1$ to orthonormal bases $\{\mathbf{u}_1, \dots, \mathbf{u}_N\}, \{\mathbf{v}_1, \dots, \mathbf{v}_N\}$ for \mathbf{R}^N , respectively. Then, the real part $\mathbf{x} \in \mathbf{R}^N$ and the imaginary

part $\mathbf{y} \in \mathbf{R}^N$ of the input signal $\mathbf{z} = \mathbf{x} + i\mathbf{y} \in \mathbf{C}^N$ can be decomposed as

$$\mathbf{x} = \sum_{i=1}^N \chi_i \mathbf{u}_i, \quad (39)$$

$$\mathbf{y} = \sum_{i=1}^N \psi_i \mathbf{v}_i. \quad (40)$$

Then, noticing that $\chi_i, \psi_i \sim N(0, 1)$, we have

$$\begin{aligned} E[\chi_i \chi_j] &= E[\mathbf{u}_i^T \mathbf{x} \mathbf{x}^T \mathbf{u}_j] = \mathbf{u}_i^T E[\mathbf{x} \mathbf{x}^T] \mathbf{u}_j \\ &= \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}, \end{aligned} \quad (41)$$

$$\begin{aligned} \mathbf{x} \mathbf{x}^T &= \chi_1^2 \mathbf{u}_1 \mathbf{u}_1^T + \sum_{i=2}^N \chi_1 \chi_i (\mathbf{u}_1 \mathbf{u}_i^T + \mathbf{u}_i \mathbf{u}_1^T) \\ &\quad + \sum_{i,j=2}^N \chi_i \chi_j \mathbf{u}_i \mathbf{u}_j^T, \end{aligned} \quad (42)$$

$$\mathbf{u}^T \cdot \mathbf{x} = \mathbf{u}^T (\chi_1 \mathbf{u}_1 + \dots + \chi_N \mathbf{u}_N) = \chi_1 u, \quad (43)$$

$$\mathbf{v}^T \cdot \mathbf{y} = \mathbf{v}^T (\psi_1 \mathbf{v}_1 + \dots + \psi_N \mathbf{v}_N) = \psi_1 v. \quad (44)$$

From Eqs. (41) – (44), we find that the first term of Eq. (26) (\mathbf{A}_{11}) is given by

$$\begin{aligned} E[(\varphi'(S))^2 \mathbf{x} \mathbf{x}^T] &= E[(\varphi'(\mathbf{u}^T \mathbf{x} - \mathbf{v}^T \mathbf{y} + c))^2 \mathbf{x} \mathbf{x}^T] \\ &= E[(\varphi'(\chi_1 u - \psi_1 v + c))^2 \chi_1^2] \mathbf{u}_1 \mathbf{u}_1^T \\ &\quad + E[(\varphi'(\chi_1 u - \psi_1 v + c))^2] \sum_{i=2}^N \mathbf{u}_i \mathbf{u}_i^T. \end{aligned} \quad (45)$$

Next, in order to calculate the second term of Eq. (26) (\mathbf{A}_{11}), decompose the real part $\mathbf{x} \in \mathbf{R}^N$ and the imaginary part $\mathbf{y} \in \mathbf{R}^N$ of the input signal $\mathbf{z} = \mathbf{x} + i\mathbf{y} \in \mathbf{C}^N$ as

$$\mathbf{x} = \sum_{i=1}^N \chi'_i \mathbf{v}_i, \quad (46)$$

$$\mathbf{y} = \sum_{i=1}^N \psi'_i \mathbf{u}_i. \quad (47)$$

Then, we have

$$E[\psi'_i \psi'_j] = \delta_{ij}, \quad (48)$$

$$\begin{aligned} \mathbf{y} \mathbf{y}^T &= (\psi'_1)^2 \mathbf{u}_1 \mathbf{u}_1^T + \sum_{i=2}^N \psi'_1 \psi'_i (\mathbf{u}_1 \mathbf{u}_i^T + \mathbf{u}_i \mathbf{u}_1^T) \\ &\quad + \sum_{i,j=2}^N \psi'_i \psi'_j \mathbf{u}_i \mathbf{u}_j^T, \end{aligned} \quad (49)$$

$$\mathbf{v}^T \cdot \mathbf{x} = \mathbf{v}^T (\chi'_1 \mathbf{v}_1 + \dots + \chi'_N \mathbf{v}_N) = \chi'_1 v, \quad (50)$$

$$\mathbf{u}^T \cdot \mathbf{y} = \mathbf{u}^T (\psi'_1 \mathbf{u}_1 + \dots + \psi'_N \mathbf{u}_N) = \psi'_1 u. \quad (51)$$

From Eqs. (48) – (51), we find that the second term of Eq.

(26) (\mathbf{A}_{11}) is given by

$$\begin{aligned} E[(\varphi'(T))^2 \mathbf{y} \mathbf{y}^T] &= E[(\varphi'(\mathbf{v}^T \mathbf{x} + \mathbf{u}^T \mathbf{y} + d))^2 \mathbf{y} \mathbf{y}^T] \\ &= E[(\varphi'(\chi'_1 v + \psi'_1 u + d))^2 (\psi'_1)^2] \mathbf{u}_1 \mathbf{u}_1^T \\ &\quad + E[(\varphi'(\chi'_1 v + \psi'_1 u + d))^2] \sum_{i=2}^N \mathbf{u}_i \mathbf{u}_i^T. \end{aligned} \quad (52)$$

Thus, from Eqs. (45) and (52),

$$\begin{aligned} \mathbf{A}_{11} &= E[(\varphi'(S))^2 \mathbf{x} \mathbf{x}^T] + E[(\varphi'(T))^2 \mathbf{y} \mathbf{y}^T] \\ &= \{E[(\varphi'(\chi_1 u - \psi_1 v + c))^2 \chi_1^2] \\ &\quad + E[(\varphi'(\chi'_1 v + \psi'_1 u + d))^2 \psi_1^2]\} \mathbf{u}_1 \mathbf{u}_1^T \\ &\quad + \{E[(\varphi'(\chi_1 u - \psi_1 v + c))^2] \\ &\quad + E[(\varphi'(\chi'_1 v + \psi'_1 u + d))^2]\} \sum_{i=2}^N \mathbf{u}_i \mathbf{u}_i^T \\ &= d_0(u, v, c, d) \cdot I \\ &\quad + \{d_2(u, v, c, d) - d_0(u, v, c, d)\} \cdot \frac{\mathbf{u} \mathbf{u}^T}{u^2}, \end{aligned} \quad (53)$$

where

$$d_0(u, v, c, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1 u - \psi_1 v + c)\}^2] + E[\{\varphi'(\chi'_1 v + \psi'_1 u + d)\}^2], \quad (54)$$

$$d_2(u, v, c, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1 u - \psi_1 v + c)\}^2 \chi_1^2] + E[\{\varphi'(\chi'_1 v + \psi'_1 u + d)\}^2 (\psi'_1)^2]. \quad (55)$$

Similarly, we have

$$\begin{aligned} \mathbf{A}_{21} &= d_{11}(u, v, c, d) \cdot \frac{\mathbf{v} \mathbf{u}^T}{uv} \\ &= \mathbf{A}_{12}^T, \end{aligned} \quad (56)$$

$$\begin{aligned} \mathbf{A}_{22} &= d_0(u, v, c, d) \cdot I \\ &\quad + \{d'_2(u, v, c, d) - d_0(u, v, c, d)\} \cdot \frac{\mathbf{v} \mathbf{v}^T}{v^2}, \end{aligned} \quad (57)$$

$$\begin{aligned} \mathbf{a}_{31} &= d_{1x}(u, v, c) \cdot \mathbf{u}_1^T \\ &= \mathbf{a}_{13}^T, \end{aligned} \quad (58)$$

$$\begin{aligned} \mathbf{a}_{32} &= -d_{1y}(u, v, c) \cdot \mathbf{v}_1^T \\ &= \mathbf{a}_{23}^T, \end{aligned} \quad (59)$$

$$\begin{aligned} \mathbf{a}_{41} &= d'_{1y}(u, v, d) \cdot \mathbf{u}_1^T \\ &= \mathbf{a}_{14}^T, \end{aligned} \quad (60)$$

$$\begin{aligned} \mathbf{a}_{42} &= d'_{1x}(u, v, d) \cdot \mathbf{v}_1^T \\ &= \mathbf{a}_{24}^T, \end{aligned} \quad (61)$$

$$\mathbf{A}_{44} = \begin{pmatrix} d_0(u, v, c) & 0 \\ 0 & d'_0(u, v, d) \end{pmatrix} \quad (62)$$

where

$$d_{11}(u, v, c, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1'v + \psi_1'u + d)\}^2 \chi_1' \psi_1'] - E[\{\varphi'(\chi_1u - \psi_1v + c)\}^2 \chi_1 \psi_1], \quad (63)$$

$$d_2'(u, v, c, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1u - \psi_1v + c)\}^2 \psi_1^2] + E[\{\varphi'(\chi_1'v + \psi_1'u + d)\}^2 (\chi_1')^2], \quad (64)$$

$$d_{1x}(u, v, c) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1u - \psi_1v + c)\}^2 \chi_1], \quad (65)$$

$$d_{1y}(u, v, c) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1u - \psi_1v + c)\}^2 \psi_1], \quad (66)$$

$$d_{1y}'(u, v, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1'v + \psi_1'u + d)\}^2 \psi_1], \quad (67)$$

$$d_{1x}'(u, v, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1'v + \psi_1'u + d)\}^2 \chi_1], \quad (68)$$

$$d_0(u, v, c) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1u - \psi_1v + c)\}^2], \quad (69)$$

$$d_0'(u, v, d) \stackrel{\text{def}}{=} E[\{\varphi'(\chi_1'v + \psi_1'u + d)\}^2]. \quad (70)$$

We compute the inverse of $A(\theta)$ (Eq. (37)) using the following formula used in [13].

Lemma 1:

$$\begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}^{-1} = \begin{pmatrix} B^{11} & B^{12} \\ B^{21} & B^{22} \end{pmatrix}, \quad (71)$$

provided

$$|B_{11}| \neq 0, |B_{22} - B_{21}B_{11}^{-1}B_{12}| \neq 0 \quad (72)$$

where

$$B^{11} = B_{11}^{-1} + B_{11}^{-1}B_{12}B_{22,1}^{-1}B_{21}B_{11}^{-1}, \quad (73)$$

$$B_{22,1} = B_{22} - B_{21}B_{11}^{-1}B_{12}, \quad (74)$$

$$B^{22} = B_{22,1}^{-1}, \quad (75)$$

$$B^{12} = (B^{21})^T = -B_{11}^{-1}B_{12}B_{22,1}^{-1}. \quad (76)$$

□

By using Lemma 1, we have

$$A(\theta)^{-1} = \begin{pmatrix} A^{11} & A^{12} \\ A^{21} & A^{22} \end{pmatrix}, \quad (77)$$

where

$$A^{11} = \begin{pmatrix} a^{11} & a^{12} \\ a^{21} & a^{22} \end{pmatrix}, \quad (78)$$

$$a^{11} = \frac{\mathbf{I}}{d_0} + \left\{ G + \left(\frac{1}{d_0} + G \right) (m_1l_1 + m_2l_2) + H(m_1l_3 + m_2l_4) \right\} \frac{\mathbf{u}\mathbf{u}^T}{u^2}, \quad (79)$$

$$a^{12} = \left\{ H + H(m_1l_1 + m_2l_2) + \left(\frac{1}{d_0} + E \right) (m_1l_3 + m_2l_4) \right\} \frac{\mathbf{u}\mathbf{v}^T}{uv}, \quad (80)$$

$$a^{21} = \left\{ H + \left(\frac{1}{d_0} + G \right) (m_3l_1 + m_4l_2) + H(m_3l_3 + m_4l_4) \right\} \frac{\mathbf{v}\mathbf{u}^T}{vu}, \quad (81)$$

$$a^{22} = \frac{\mathbf{I}}{d_0} + \{ E + H(m_3l_1 + m_4l_2) + \left(\frac{1}{d_0} + E \right) (m_3l_3 + m_4l_4) \} \frac{\mathbf{v}\mathbf{v}^T}{v^2}, \quad (82)$$

$$A^{12} = \begin{pmatrix} l_1\mathbf{u}/u & l_2\mathbf{u}/u \\ l_3\mathbf{v}/v & l_4\mathbf{v}/v \end{pmatrix} = (A^{21})^T, \quad (83)$$

$$A^{22} = \frac{1}{k_1k_2 - k^2} \begin{pmatrix} k_2 & -k \\ -k & k_1 \end{pmatrix}, \quad (84)$$

$$E = \frac{d_2}{d_2d_2' - d_{11}^2} - \frac{1}{d_0}, \quad (85)$$

$$F = \frac{1}{d_2} - \frac{1}{d_0}, \quad (86)$$

$$G = Ed_{11}^2 \left(\frac{1 + 2F}{d_0} + F^2 \right) + \frac{d_{11}^2}{d_0} \left(F^2 + 2F + \frac{1}{d_0^2} \right) + F, \quad (87)$$

$$H = -d_{11} \left(\frac{1}{d_0} + E \right) \left(\frac{1}{d_0} + F \right), \quad (88)$$

$$k = -d_{1x}d_{1y}' \left(\frac{1}{d_0} + G \right) - H(d_{1y}d_{1y}' - d_{1x}d_{1x}') + d_{1x}'d_{1y} \left(\frac{1}{d_0} + E \right), \quad (89)$$

$$k_1 = d_0 - d_{1x}^2 \left(\frac{1}{d_0} + G \right) + 2Hd_{1x}d_{1y} - d_{1y}^2 \left(\frac{1}{d_0} + E \right), \quad (90)$$

$$k_2 = d_0' - d_{1y}'^2 \left(\frac{1}{d_0} + G \right) - 2Hd_{1x}'d_{1y}' - d_{1x}'^2 \left(\frac{1}{d_0} + E \right), \quad (91)$$

$$l_1 = \frac{1}{k_1k_2 - k^2} \left\{ -k_2 \left(\frac{d_{1x}}{d_0} + Gd_{1x} - Hd_{1y} \right) + k \left(\frac{d_{1y}'}{d_0} + Gd_{1y}' + Hd_{1x}' \right) \right\}, \quad (92)$$

ACKNOWLEDGMENT

The author would like to give special thanks to the anonymous reviewers for valuable comments.

REFERENCES

- [1] S. Amari, *Differential-geometrical methods in statistics*, Lecture Notes in Statistics vol. 28, Springer-Verlag, 1985.
- [2] S. Amari, "Neural learning in structured parameter spaces – Natural Riemannian gradient," In M. C. Mozer, M. I. Jordan, & Th. Petsche (Eds.), *Advances in neural processing systems*, 9, Cambridge, MA: MIT Press, 1996.
- [3] S. Amari, "Natural gradient works efficiently in learning," *Neural Computation*, vol. 10, no. 2, pp. 251–276, 1998.
- [4] S. Amari, H. Park, and T. Ozeki, "Singularities affect dynamics of learning in neuromanifolds," *Neural Computation*, vol. 18, no. 5, pp. 1007–1065, 2006.
- [5] G. M. Georgiou and C. Koutsougeras, "Complex domain backpropagation," *IEEE Trans. Circuits and Systems-II: Analog and Digital Signal Processing*, vol. 39, no. 5, pp. 330–334, 1992.
- [6] T. Nitta (Ed.), *Complex-valued neural networks: utilizing high-dimensional parameters*, Information Science Reference, Pennsylvania, USA, 2009.
- [7] T. Nitta, "Local minima in hierarchical structures of complex-valued neural networks," *Neural Networks*, vol. 43, pp. 1–7, 2013.
- [8] T. Nitta, "Plateau in a polar variable complex-valued neuron," Proceedings of the 6th International Conference on Agents and Artificial Intelligence, ICAART2014-Anger, March 6-8, 2014, pp. 526–531.
- [9] A. Hirose (Ed.), *Complex-valued neural networks: advances and applications in the IEEE press series on computational intelligence*, Wiley-IEEE Press, 2013.
- [10] M. S. Kim and C. C. Guest, "Modification of backpropagation networks for complex-valued signal processing in frequency domain," *Proc. International Joint Conference on Neural Networks*, San Diego, June, 1990, vol. 3, pp. 27–31.
- [11] T. Nitta, "An extension of the back-propagation algorithm to complex numbers," *Neural Networks*, vol. 10, no. 8, pp. 1392–1415, 1997.
- [12] T. Nitta, "Orthogonality of decision boundaries in complex-valued neural networks," *Neural Computation*, vol. 16, no. 1, pp. 73–97, 2004.
- [13] H. H. Yang and S. Amari, "Complexity issues in natural gradient descent method for training multilayer perceptrons," *Neural Computation*, vol. 10, no. 8, pp. 2137–2157, 1998.

V. CONCLUDING REMARK

We have developed the natural gradient descent method for the multilayer stochastic complex-valued neural networks, and derived the natural gradient for a single stochastic complex-valued neuron. Since we assume that the variance σ_1^2 for the real part of the complex-valued additive noise is equal to that σ_2^2 for the imaginary, the Riemannian metric tensor is given by $\mathbf{G}(\theta) = (1/\sigma^2)\mathbf{A}(\theta)$ where $\mathbf{A}(\theta)$ does not depend on $\sigma^2 = \sigma_1^2 = \sigma_2^2$. The situation is, however, not the same if $\sigma_1 \neq \sigma_2$. And also, the Riemannian metric tensor will be more complicated if the covariance is not zero.

In future studies, based on the results of this paper, we will derive the natural gradient descent algorithm for the three-layered stochastic complex-valued neural network, and make clear its properties via computer simulations. Since there exist many singular points in the three-layered stochastic complex-valued neural network [7], it is expected that the complex-valued natural gradient method improves its learning performance dramatically. And also, it has been shown that there exist singular points in the polar variable complex-valued neurons [8]. We will apply the natural gradient method to the polar variable complex-valued neurons.