

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2577260>

# Natural Gradient Descent for Training Multi-Layer Perceptrons

Article · January 1998

Source: CiteSeer

---

CITATIONS

18

---

READS

304

3 authors, including:



**Howard H Yang**

National Institutes of Health

240 PUBLICATIONS 6,924 CITATIONS

SEE PROFILE



**Shun-ichi Amari**

RIKEN

564 PUBLICATIONS 34,909 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Information Geometry and Applications to Neuroimaging [View project](#)



state-space Ising model [View project](#)

# Natural Gradient Descent for Training Multi-Layer Perceptrons

Howard Hua Yang and Shun-ichi Amari  
*Lab. for Information Representation, FRP, RIKEN*  
*Hirosawa 2-1, Wako-shi, Saitama 351-01, JAPAN*  
FAX: +81 48462 4633  
E-mails: hhy@koala.riken.go.jp, amari@zoo.riken.go.jp

Submitted to IEEE Tr. on Neural Networks  
May 16, 1996

## Abstract

The main difficulty in implementing the natural gradient learning rule is to compute the inverse of the Fisher information matrix when the input dimension is large. We have found a new scheme to represent the Fisher information matrix. Based on this scheme, we have designed an algorithm to compute the inverse of the Fisher information matrix. When the input dimension  $n$  is much larger than the number of hidden neurons, the complexity of this algorithm is of order  $O(n^2)$  while the complexity of conventional algorithms for the same purpose is of order  $O(n^3)$ . The simulation has confirmed the efficiency and robustness of the natural gradient learning rule.

## 1 Introduction

Inversion of the Fisher information matrix is required to obtain the Cramer-Rao lower bound which is fundamental for analyzing the performance of an unbiased estimator. It is also needed in the natural gradient learning framework [2, 3] to design statistically efficient learning algorithms for parameter estimation in general and for training neural networks in particular. In this paper, we assume a stochastic model for multi-layer perceptrons. Considering the parameter space as a Riemannian space with the Fisher information matrix as its metric, we apply the natural gradient learning rule to train a multi-layer perceptron. The main difficulty encountered is to compute the inverse of the Fisher information matrix of a large size when the input dimension is large. By exploring the structure of the Fisher information matrix and its inverse, we shall design a fast algorithm with lower complexity to implement the natural gradient learning algorithm.

The outline of the paper is as follows. The natural gradient descent method and related problems for training a multi-layer perceptron are described in Section 2. The Fisher information matrix and the scheme to represent this matrix are given in Section 3. The main focus of this paper is to present a fast algorithm to compute the inverse of the Fisher information matrix. Several issues related to this are discussed in Section 4: a group structure for the blocks in the Fisher information matrix is found in Section 4.1. Based on this finding the inverse of the Fisher information matrix is computed in Section 4.2-4.3; the complexity of our algorithm in computing the inverse of the Fisher information matrix is analyzed in Section 4.4; by the

result on the group structure of the Fisher information matrix, the Cramer-Rao lower bound for the committee machine is explicitly given in Section 4.5. The natural gradient vector field is compared with the ordinary vector field in Section 5 via an example to illustrate the effects of the natural gradient learning in training multi-layer perceptrons. The effectiveness of the natural gradient is further demonstrated by some simulation results in Section 6.

## 2 Natural Gradient Descent Method and Its Efficiency

### 2.1 Model of a stochastic multi-layer perceptron

Assume the following model of a stochastic multi-layer perceptron:

$$z = \sum_{i=1}^m a_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i) + \xi \quad (1)$$

where  $(\cdot)^T$  denotes the transpose,  $\xi \sim N(0, \sigma^2)$  is a Gaussian random variable, and  $\varphi(x)$  is a differentiable output function of hidden neurons such as hyperbolic tangent. Assume that the multi-layer network has a  $n$ -dimensional input,  $m$  hidden neurons, a one dimensional output, and  $m \leq n$ . Denote  $\mathbf{a} = (a_1, \dots, a_m)^T$  the weight vector of the output neuron,  $\mathbf{w}_i = (w_{1i}, \dots, w_{ni})^T$  the weight vector of the  $i$ -th hidden neuron, and  $\mathbf{b} = (b_1, \dots, b_m)^T$  is the vector of thresholds for the hidden neurons. Let  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_m]$  be a matrix formed by column weight vectors  $\mathbf{w}_i$ , then (1) can be rewritten as

$$z = \mathbf{a}^T \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) + \xi.$$

Here, the scalar function  $\varphi$  operates on each component of the vector  $\mathbf{W}^T \mathbf{x} + \mathbf{b}$  and this convention will be used throughout this paper. By this convention, the column vector

$$[\varphi(\mathbf{w}_1^T \mathbf{x} + b_1), \dots, \varphi(\mathbf{w}_m^T \mathbf{x} + b_m)]^T$$

is written as  $\varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b})$  and the row vector

$$[\varphi(\mathbf{w}_1^T \mathbf{x} + b_1), \dots, \varphi(\mathbf{w}_m^T \mathbf{x} + b_m)]$$

is written as  $\varphi(\mathbf{x}^T \mathbf{W} + \mathbf{b}^T)$ .

In this paper, we consider the realizable learning problem. Let

$$\boldsymbol{\theta} = (\mathbf{w}_1^T, \dots, \mathbf{w}_m^T, \mathbf{a}^T, \mathbf{b}^T)^T$$

denote the parameter vector for the student network and

$$\boldsymbol{\theta}^* = (\mathbf{w}_1^{*T}, \dots, \mathbf{w}_m^{*T}, \mathbf{a}^{*T}, \mathbf{b}^{*T})^T$$

for the teacher network.

The problem is to estimate  $\boldsymbol{\theta}^*$  based on training examples

$$D_T = \{(\mathbf{x}_t, z_t), t = 1, \dots, T\}.$$

The training examples as well as the future examples for evaluating the generalization error are generated by two steps. First,  $\mathbf{x}_t$  is randomly generated subject to an unknown probability

density function  $p(\mathbf{x})$ . Second, the corresponding output is generated by the teacher network as

$$z_t = \mathbf{a}^{*T} \varphi(\mathbf{W}^{*T} \mathbf{x}_t + \mathbf{b}^*) + \xi_t$$

where  $\xi_t \sim N(0, \sigma^2)$ . Note the variance of the additive noise  $\xi_t$  is the same as the variance of the additive noise  $\xi$  in the stochastic perceptron model.

The input-output pairs  $\{(\mathbf{x}_t, z_t), t > T\}$  are the future examples. Define the generalization error

$$E_g(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = E_{\boldsymbol{\theta}^*} [(z - \mathbf{a}^T \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}))^2]$$

where the expectation  $E_{\boldsymbol{\theta}^*}[\cdot]$  is taken with respect to the joint probability density function  $p(\mathbf{x}, z; \boldsymbol{\theta}^*)$ .

The purpose of learning is to minimize the generalization error but it cannot be easily evaluated in practice. Instead of the generalization error, we use the training error

$$E_{\text{tr}}(\boldsymbol{\theta}, D_T) = \frac{1}{T} \sum_{t=1}^T \{z_t - \mathbf{a}^T \varphi(\mathbf{W}^T \mathbf{x}_t + \mathbf{b})\}^2$$

to derive a learning algorithm. Minimizing the training error is not equivalent to minimizing the generalization error in general (see, for example, Amari and Murata, 1993)[4]. Some auxiliary methods such as early stopping and regularization are necessary to avoid overtraining or overfitting (see, for example, Ripley, 1996 [10]; Amari et al, 1997 [5]). The overtraining problem is important but it will not be discussed here. We shall only focus the learning algorithm for minimizing the training error.

## 2.2 Fisher information matrix

The joint probability density function (pdf) of the input and the output is

$$p(\mathbf{x}, z; \mathbf{W}, \mathbf{a}, \mathbf{b}) = p(z|\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b})p(\mathbf{x})$$

where  $p(z|\mathbf{x}; \mathbf{W}, \mathbf{a}, \mathbf{b})$  is the conditional pdf of  $z$  when  $\mathbf{x}$  is input. The loss function is defined as the negative log-likelihood function

$$L(\mathbf{x}, z; \boldsymbol{\theta}) = -\log p(\mathbf{x}, z; \boldsymbol{\theta}) = l(z|\mathbf{x}; \boldsymbol{\theta}) - \log p(\mathbf{x})$$

where

$$l(z|\mathbf{x}; \boldsymbol{\theta}) = -\log p(z|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{2\sigma^2} (z - \mathbf{a}^T \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}))^2$$

is also a loss function equivalent to  $L(\mathbf{x}, z; \boldsymbol{\theta})$  since  $p(\mathbf{x})$  does not depend on  $\boldsymbol{\theta}$ . Hence, given the training set  $D_T$ , minimizing the loss function  $L(\mathbf{x}, z; \boldsymbol{\theta})$  is equivalent to minimizing the training error.

Since  $\frac{\partial L}{\partial \boldsymbol{\theta}} = \frac{\partial l}{\partial \boldsymbol{\theta}}$ , the Fisher information matrix is defined by

$$\mathbf{G} = \mathbf{G}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}} \left[ \frac{\partial L}{\partial \boldsymbol{\theta}} \left( \frac{\partial L}{\partial \boldsymbol{\theta}} \right)^T \right] = E_{\boldsymbol{\theta}} \left[ \frac{\partial l}{\partial \boldsymbol{\theta}} \left( \frac{\partial l}{\partial \boldsymbol{\theta}} \right)^T \right] \quad (2)$$

where  $E_{\boldsymbol{\theta}}[\cdot]$  denotes the expectation with respect to  $p(\mathbf{x}, z; \boldsymbol{\theta})$ .

The inverse of the Fisher information matrix is often used in performance analysis, because the Cramer-Rao states that an unbiased estimator  $\hat{\boldsymbol{\theta}}$  of the true parameter  $\boldsymbol{\theta}^*$  satisfies

$$E_{\boldsymbol{\theta}^*} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|^2] \geq \text{Tr}(\mathbf{G}^{-1}(\boldsymbol{\theta}^*))$$

where  $\text{Tr}(\cdot)$  denotes the trace of a matrix.

For the on-line estimator  $\hat{\boldsymbol{\theta}}_t$  based on the independent examples  $\{(\mathbf{x}_s, z_s), s = 1, \dots, t\}$ , the negative log-likelihood function is

$$L_{1,t} = L(\mathbf{x}_1, \dots, \mathbf{x}_t, z_1, \dots, z_t; \boldsymbol{\theta}) = -\log p(\mathbf{x}_1, \dots, \mathbf{x}_t, z_1, \dots, z_t; \boldsymbol{\theta}) = \sum_{s=1}^t L(\mathbf{x}_s, z_s; \boldsymbol{\theta}).$$

Since

$$E_{\boldsymbol{\theta}} \left[ \frac{\partial L}{\partial \boldsymbol{\theta}}(\mathbf{x}_s, z_s; \boldsymbol{\theta}) \left( \frac{\partial L}{\partial \boldsymbol{\theta}}(\mathbf{x}_\tau, z_\tau; \boldsymbol{\theta}) \right)^T \right] = \delta_{s,\tau} \mathbf{G}(\boldsymbol{\theta})$$

due to the independence of the examples,

$$\mathbf{G}_t(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}} \left[ \frac{\partial L_{1,t}}{\partial \boldsymbol{\theta}} \left( \frac{\partial L_{1,t}}{\partial \boldsymbol{\theta}} \right)^T \right] = t \mathbf{G}(\boldsymbol{\theta}).$$

So when the examples  $\{(\mathbf{x}_s, z_s), s = 1, 2, \dots\}$  are independently drawn according to the probability law  $p(\mathbf{x}, z; \boldsymbol{\theta}^*)$ , the Cramer-Rao inequality for the on-line estimator is

$$E_{\boldsymbol{\theta}^*} [\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|^2] \geq \frac{1}{t} \text{Tr}(\mathbf{G}^{-1}(\boldsymbol{\theta}^*)) \quad (3)$$

### 2.3 Natural gradient learning

Let  $\Theta = \{\boldsymbol{\theta}\}$  be the parameter space. The divergence between two points  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\theta}_2$ , or two distributions,  $p(\mathbf{x}, z; \boldsymbol{\theta}_1)$  and  $p(\mathbf{x}, z; \boldsymbol{\theta}_2)$ , is given by the Kullback-Leibler divergence

$$D(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \text{KL}[p(\mathbf{x}, z; \boldsymbol{\theta}_1) \| p(\mathbf{x}, z; \boldsymbol{\theta}_2)].$$

When the two points are infinitesimally close, we have the quadratic form

$$D(\boldsymbol{\theta}, \boldsymbol{\theta} + d\boldsymbol{\theta}) = \frac{1}{2} d\boldsymbol{\theta}^T \mathbf{G}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (4)$$

where  $\mathbf{G}(\boldsymbol{\theta})$  is the Fisher information matrix. This is regarded as the square of the length of  $d\boldsymbol{\theta}$ . Since  $\mathbf{G}(\boldsymbol{\theta})$  depends on  $\boldsymbol{\theta}$ , the parameter space is regarded as a Riemannian space in which the local distance is defined by (4). Here, the Fisher information matrix  $\mathbf{G}(\boldsymbol{\theta})$  plays the role of the Riemannian metric tensor. See (Amari, 1985)[1] and (Murray and Rice, 1993)[9] for information geometry.

The natural gradient descent method proposed by Amari [3, 2] makes use of the Riemannian metric given by the Fisher information matrix to optimize the learning dynamics such that the Cramer-Rao lower bound is achieved asymptotically. The idea is to convert the covariant gradient  $\frac{\partial l}{\partial \boldsymbol{\theta}}$  into contravariant form  $\mathbf{G}^{-1} \frac{\partial l}{\partial \boldsymbol{\theta}}$ . It is shown by Amari [2] that the steepest descent direction of a function  $C(\boldsymbol{\theta})$  in the Riemannian space  $\Theta$  is

$$-\tilde{\nabla} C(\boldsymbol{\theta}) = -\mathbf{G}^{-1}(\boldsymbol{\theta}) \nabla C(\boldsymbol{\theta}).$$

The on-line learning algorithms corresponding to  $\frac{\partial l}{\partial \boldsymbol{\theta}}$  and  $\mathbf{G}^{-1}(\boldsymbol{\theta}) \frac{\partial l}{\partial \boldsymbol{\theta}}$  are, respectively, the ordinary gradient descent algorithm:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\mu}{t} \frac{\partial l}{\partial \boldsymbol{\theta}}(z_t | \mathbf{x}_t; \boldsymbol{\theta}_t) \quad (5)$$

and the natural gradient descent algorithm:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\mu}{t} \mathbf{G}^{-1}(\boldsymbol{\theta}_t) \frac{\partial l}{\partial \boldsymbol{\theta}}(z_t | \mathbf{x}_t; \boldsymbol{\theta}_t) \quad (6)$$

where  $\mu$  and  $\mu'$  are learning rates.

It is proved in [2] that the natural gradient learning is Fisher efficient. When the negative log-likelihood function is taken as the loss function, the natural gradient descent algorithm (6) gives an efficient on-line estimator, i.e., the asymptotic variance of  $\boldsymbol{\theta}_t$  driven by (6) satisfies

$$E_{\boldsymbol{\theta}^*}[(\boldsymbol{\theta}_t - \boldsymbol{\theta}^*)(\boldsymbol{\theta}_t - \boldsymbol{\theta}^*)^T] \approx \frac{1}{t} \mathbf{G}^{-1}(\boldsymbol{\theta}^*) \quad (7)$$

which gives the mean square error

$$E_{\boldsymbol{\theta}^*}[\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|^2] \approx \frac{1}{t} \text{Tr}(\mathbf{G}^{-1}(\boldsymbol{\theta}^*)). \quad (8)$$

However, the algorithm (5) does not attain this limit.

**Remark 1** *The equilibrium points are the same for the averaged versions of the dynamics (5) and (6), but the inverse of the Fisher information matrix optimizes the learning dynamics. The local minima problem still exists in the natural gradient learning.*

**Remark 2** *The natural gradient learning algorithm is efficient, but its complexity is generally high due to the computation of the matrix inverse. In later sections, we shall explore the structure of the Fisher information matrix for a multi-layer perceptron and design a low complexity algorithm to compute  $\mathbf{G}^{-1}(\boldsymbol{\theta})$ . The dimension of  $\mathbf{G}(\boldsymbol{\theta})$  is  $(nm + 2m) \times (nm + 2m)$ , the time complexity of computing  $\mathbf{G}^{-1}(\boldsymbol{\theta})$  by generally accepted algorithms is  $O(n^3)$  which is too high for the on-line natural gradient learning algorithm. Our method reduces the time complexity to  $O(n^2)$ .*

### 3 Fisher Information Matrix of A Multi-Layer Perceptron

#### 3.1 Likelihood function of the stochastic multi-layer perceptron

In previous sections, we did not give any specific assumption about the input. To compute an explicit expression of the Fisher information matrix, we now assume that the input  $\mathbf{x}$  is subject to the standard Gaussian distribution  $N(\mathbf{0}, \mathbf{I})$ , where  $\mathbf{I}$  is the identity matrix.

From  $l(z|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{2\sigma^2}(z - \sum_{i=1}^m a_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i))^2$ , we obtain

$$\frac{\partial l}{\partial w_{ij}} = -\frac{1}{\sigma^2} \left( z - \sum_{k=1}^m a_k \varphi(\mathbf{w}_k^T \mathbf{x} + b_k) \right) a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) x_j \quad (9)$$

$$\frac{\partial l}{\partial a_i} = -\frac{1}{\sigma^2} \left( z - \sum_{k=1}^m a_k \varphi(\mathbf{w}_k^T \mathbf{x} + b_k) \right) \varphi(\mathbf{w}_i^T \mathbf{x} + b_i) \quad (10)$$

$$\frac{\partial l}{\partial b_i} = -\frac{1}{\sigma^2} \left( z - \sum_{k=1}^m a_k \varphi(\mathbf{w}_k^T \mathbf{x} + b_k) \right) a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \quad (11)$$

When  $(\mathbf{x}, z)$  is the input-output pair of the stochastic multi-layer perceptron, the equations (9)-(11) become

$$\frac{\partial l}{\partial w_{ij}} = -\frac{\xi}{\sigma^2} a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) x_j \quad (12)$$

$$\frac{\partial l}{\partial a_i} = -\frac{\xi}{\sigma^2} \varphi(\mathbf{w}_i^T \mathbf{x} + b_i) \quad (13)$$

$$\frac{\partial l}{\partial b_i} = -\frac{\xi}{\sigma^2} a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i), \quad (14)$$

and in vector forms

$$\frac{\partial l}{\partial \mathbf{w}_i} = -\frac{\xi}{\sigma^2} a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \mathbf{x}, i = 1, \dots, m \quad (15)$$

$$\frac{\partial l}{\partial \mathbf{a}} = -\frac{\xi}{\sigma^2} \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \quad (16)$$

$$\frac{\partial l}{\partial \mathbf{b}} = -\frac{\xi}{\sigma^2} \{\mathbf{a} \odot \varphi'(\mathbf{W}^T \mathbf{x} + \mathbf{b})\} \quad (17)$$

where  $\odot$  denotes the Hadmard product which is the componentwise product of vectors or matrices. The above three vector forms can be written as one:

$$\frac{\partial l}{\partial \boldsymbol{\theta}} = -\frac{\xi}{\sigma^2} \begin{bmatrix} (\mathbf{a} \odot \varphi'(\mathbf{W}^T \mathbf{x} + \mathbf{b})) \otimes \mathbf{x} \\ \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \\ \mathbf{a} \odot \varphi'(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \end{bmatrix} \quad (18)$$

where  $\otimes$  denotes the Kronecker product. For two matrices  $\mathbf{A} = (a_{ij})$  and  $\mathbf{B} = (b_{ij})$ ,  $\mathbf{A} \otimes \mathbf{B}$  is defined as a larger matrix consisting of the block matrices  $[a_{ij} \mathbf{B}]$ . For two column vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,  $\mathbf{a} \otimes \mathbf{b}$  gives the higher-dimensional column vector  $[a_1 \mathbf{b}^T, \dots, a_n \mathbf{b}^T]^T$ .

### 3.2 Blocks in the Fisher information matrix

The parameter  $\boldsymbol{\theta}$  consists of  $m + 2$  blocks

$$\boldsymbol{\theta} = [\mathbf{w}_1^T, \dots, \mathbf{w}_m^T, \mathbf{a}^T, \mathbf{b}^T]^T$$

where the first  $m$  blocks  $\mathbf{w}_i$ 's are  $n$ -dimensional and the last two blocks  $\mathbf{a}$  and  $\mathbf{b}$  are  $m$ -dimensional. From (18), the Fisher information matrix  $\mathbf{G}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \mathbf{A}(\boldsymbol{\theta})$  is of  $(n+2)m \times (n+2)m$  dimensions. Here, the matrix  $\mathbf{A}(\boldsymbol{\theta})$  does not depend on  $\sigma^2$ . It is partitioned into  $(m+2) \times (m+2)$  blocks of sub-matrices corresponding to the partition of  $\boldsymbol{\theta}$ . We denote each sub-matrix by  $\mathbf{A}_{ij}$ . The Fisher information matrix is written as

$$\mathbf{G}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} [\mathbf{A}_{ij}]_{(m+2) \times (m+2)}. \quad (19)$$

The partition diagram for the matrix  $\mathbf{A}(\boldsymbol{\theta})$  is given below:

$$\mathbf{A}(\boldsymbol{\theta}) = \begin{array}{c|cccccc} & \mathbf{w}_1^T & \cdots & \mathbf{w}_m^T & \mathbf{a}^T & \mathbf{b}^T \\ \hline \mathbf{w}_1 & \mathbf{A}_{1,1} & \cdots & \mathbf{A}_{1,m} & \mathbf{A}_{1,m+1} & \mathbf{A}_{1,m+2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{w}_m & \mathbf{A}_{m,1} & \cdots & \mathbf{A}_{m,m} & \mathbf{A}_{m,m+1} & \mathbf{A}_{m,m+2} \\ \mathbf{a} & \mathbf{A}_{m+1,1} & \cdots & \mathbf{A}_{m+1,m} & \mathbf{A}_{m+1,m+1} & \mathbf{A}_{m+1,m+2} \\ \mathbf{b} & \mathbf{A}_{m+2,1} & \cdots & \mathbf{A}_{m+2,m} & \mathbf{A}_{m+2,m+1} & \mathbf{A}_{m+2,m+2} \end{array} \quad (20)$$

where

$$\begin{aligned} \mathbf{A}_{ij} &= a_i a_j E[\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \varphi'(\mathbf{w}_j^T \mathbf{x} + b_j) \mathbf{x} \mathbf{x}^T], \quad \text{for } 1 \leq i, j \leq m, \\ \mathbf{A}_{i,m+1} &= a_i E[\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \mathbf{x} \varphi(\mathbf{x}^T \mathbf{W} + \mathbf{b})] \text{ and} \\ \mathbf{A}_{m+1,i} &= \mathbf{A}_{i,m+1}^T \quad \text{for } 1 \leq i \leq m, \\ \mathbf{A}_{m+1,m+1} &= E[\varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \varphi(\mathbf{x}^T \mathbf{W} + \mathbf{b})], \\ \mathbf{A}_{i,m+2} &= E[a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \mathbf{x} (\mathbf{a}^T \odot \varphi'(\mathbf{x}^T \mathbf{W} + \mathbf{b}^T))] \text{ and} \\ \mathbf{A}_{m+2,i} &= \mathbf{A}_{i,m+2}^T \quad \text{for } 1 \leq i \leq m, \\ \mathbf{A}_{m+1,m+2} &= \mathbf{A}_{m+2,m+1}^T = E[\varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) (\mathbf{a}^T \odot \varphi'(\mathbf{x}^T \mathbf{W} + \mathbf{b}^T))], \\ \mathbf{A}_{m+2,m+2} &= (\mathbf{a} \mathbf{a}^T) \odot E[\varphi'(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \varphi'(\mathbf{x}^T \mathbf{W} + \mathbf{b}^T)]. \end{aligned}$$

Here,  $E[\cdot]$  denotes the expectation with respect to  $p(\mathbf{x})$ . From the above expressions, we know that  $\mathbf{A}(\boldsymbol{\theta}) = [\mathbf{A}_{ij}]$  does not depend on  $\sigma^2$ . The explicit forms of these matrix blocks in the Fisher information matrix are characterized by the three lemmas in the next section.

**Remark 3** *Since the loss function  $l(z|\mathbf{x};\boldsymbol{\theta})$  contains the unknown parameter  $\sigma^2$ , we should use the following modified loss function in practice:*

$$l_1(z|\mathbf{x};\boldsymbol{\theta}) = \frac{1}{2}(z - \mathbf{a}^T \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}))^2.$$

Modifying the learning rule (5), we have the following on-line algorithm:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\mu}{t} \frac{\partial l_1}{\partial \boldsymbol{\theta}}(z_t|\mathbf{x}_t; \boldsymbol{\theta}_t). \quad (21)$$

Since  $\mathbf{G}^{-1}(\boldsymbol{\theta}) \frac{\partial l}{\partial \boldsymbol{\theta}} = \mathbf{A}^{-1}(\boldsymbol{\theta}) \frac{\partial l_1}{\partial \boldsymbol{\theta}}$ , the learning rule (6) is exactly the same as the following learning rule:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\mu'}{t} \mathbf{A}^{-1}(\boldsymbol{\theta}_t) \frac{\partial l_1}{\partial \boldsymbol{\theta}}(z_t|\mathbf{x}_t; \boldsymbol{\theta}_t). \quad (22)$$

### 3.3 Representation of blocks in the Fisher information matrix

**Lemma 1** *For  $i = 1, \dots, m$ , the diagonal blocks are given by*

$$\mathbf{A}_{ii} = a_i^2 [d_1(w_i, b_i) \mathbf{I} + \{d_2(w_i, b_i) - d_1(w_i, b_i)\} \mathbf{u}_i \mathbf{u}_i^T] \quad (23)$$

where

$$w_i = \|\mathbf{w}_i\| \text{ denoting the Euclidean norm of } \mathbf{w}_i, \mathbf{u}_i = \mathbf{w}_i/w_i, \quad (24)$$

$$d_1(w, b) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\varphi'(wx + b))^2 e^{-\frac{x^2}{2}} dx > 0,$$

$$d_2(w, b) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\varphi'(wx + b))^2 x^2 e^{-\frac{x^2}{2}} dx > 0. \quad (25)$$

When  $a_i \neq 0$ , the inverse of the matrix  $\mathbf{A}_{ii}$  is given by

$$\mathbf{A}_{ii}^{-1} = \frac{1}{a_i^2} \left[ \frac{1}{d_1(w_i, b_i)} \mathbf{I} + \left\{ \frac{1}{d_2(w_i, b_i)} - \frac{1}{d_1(w_i, b_i)} \right\} \mathbf{u}_i \mathbf{u}_i^T \right], \quad i = 1, \dots, m. \quad (26)$$

The proof of Lemma 1 is given in Appendix 1.

In particular, for a single-layer stochastic perceptron

$$z = \varphi(\mathbf{w}^T \mathbf{x} + b) + \xi,$$

the expressions (23) and (26) give the Fisher information matrix:

$$\mathbf{G}(\mathbf{w}) = \frac{1}{\sigma^2} [d_1(w, b) \mathbf{I} + \{d_2(w, b) - d_1(w, b)\} \mathbf{u}_i \mathbf{u}_i^T] \quad (27)$$

and its inverse formula

$$\mathbf{G}^{-1}(\mathbf{w}) = \sigma^2 \left[ \frac{1}{d_1(w, b)} \mathbf{I} + \left\{ \frac{1}{d_2(w, b)} - \frac{1}{d_1(w, b)} \right\} \mathbf{u}_i \mathbf{u}_i^T \right]. \quad (28)$$

When  $b = 0$ , the above formula is found by Amari in [3].

To obtain the explicit forms of the other blocks  $\mathbf{A}_{ij}$  in  $\mathbf{G}$ , we need to introduce two bases in  $\mathfrak{R}^n$  which possess certain properties. Except for a set in  $\mathfrak{R}^n$  with zero Lebesgue measure, the vectors  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$  are linearly independent. It is easy to supplement  $n - m$  vectors  $\mathbf{u}_{m+1}, \dots, \mathbf{u}_n$  to them, such that they together form a basis in  $\mathfrak{R}^n$ ,  $\{\mathbf{u}_1, \dots, \mathbf{u}_m; \mathbf{u}_{m+1}, \dots, \mathbf{u}_n\}$ , having the following properties:

$$\text{for } j > m, \quad \mathbf{u}_j - \mathcal{L}(\mathbf{u}_1, \dots, \mathbf{u}_m) \quad (29)$$

$$\mathbf{u}_j^T \mathbf{u}_k = \delta_{j,k} \text{ (delta notation),} \quad \text{for } j, k = m + 1, \dots, n, \quad (30)$$

where  $\mathcal{L}(\mathbf{u}_1, \dots, \mathbf{u}_m)$  is the vector space spanned by  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ .

Let  $\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_n]$  and  $\mathbf{V} = (\mathbf{U}^{-1})^T = [\mathbf{v}_1 \cdots \mathbf{v}_n]$ . The identity

$$\mathbf{V}^T \mathbf{U} = \mathbf{U}^T \mathbf{V} = \mathbf{I}$$

implies

$$\mathbf{v}_j^T \mathbf{u}_k = \mathbf{u}_j^T \mathbf{v}_k = \delta_{j,k}. \quad (31)$$

Hence,  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is the orthogonal conjugate basis of  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ .

From (29), (30) and (31), we can easily prove that

$$\mathbf{u}_j = \mathbf{v}_j \quad \text{for } j > m.$$

The random input  $\mathbf{x}$  is represented in dual ways by

$$\mathbf{x} = \sum_{i=1}^n x_i \mathbf{u}_i = \sum_{i=1}^n x'_i \mathbf{v}_i$$

where  $(x_1, \dots, x_n)$  and  $(x'_1, \dots, x'_n)$  are coordinates of  $\mathbf{x}$  on the bases  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$  and  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  respectively. It follows from (31) that

$$x_i = \mathbf{v}_i^T \mathbf{x} = \mathbf{x}^T \mathbf{v}_i \quad \text{and} \quad x'_i = \mathbf{u}_i^T \mathbf{x} = \mathbf{x}^T \mathbf{u}_i. \quad (32)$$

In this paper, a  $m \times n$  matrix is denoted by  $(a_{ij})_{m \times n}$ . The subscript  $m \times n$  may be omitted when the dimension of the matrix can be determined by the context. Define a  $n \times n$  matrix  $\mathbf{R} = (r_{ij}) = (\mathbf{u}_i^T \mathbf{u}_j)$ . Let  $\mathbf{R}^{-1} = (r^{ij})$ . Noticing that  $\mathbf{R} = (\mathbf{u}_i^T \mathbf{u}_j) = \mathbf{U}^T \mathbf{U}$  and  $\mathbf{V} = (\mathbf{U}^{-1})^T$ , we have  $\mathbf{R}^{-1} = \mathbf{U}^{-1} (\mathbf{U}^{-1})^T = \mathbf{V}^T \mathbf{V}$ , i.e.,  $r^{ij} = \mathbf{v}_i^T \mathbf{v}_j$ . Due to the properties (29) and (30),

$$\mathbf{R} = \begin{bmatrix} (r_{ij})_{m \times m} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}_{n \times n}, \quad \text{and} \quad \mathbf{R}^{-1} = \begin{bmatrix} (r^{ij})_{m \times m} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}_{n \times n}$$

where  $\mathbf{0}$  and  $\mathbf{I}$  are zero and identity matrices respectively with proper dimensions admissible in the above partitions of  $\mathbf{R}$  and  $\mathbf{R}^{-1}$ . By the definition of  $\mathbf{R}$ , we have  $x'_j = \mathbf{u}_j^T \mathbf{x} = \mathbf{u}_j^T \sum_{k=1}^n x_k \mathbf{u}_k = \sum_{k=1}^n r_{jk} x_k$  and  $x_j = \sum_{k=1}^n r^{jk} x'_k$ . Therefore, for  $1 \leq j \leq m$ ,  $x'_j = \sum_{k=1}^m r_{jk} x_k$  and  $x_j = \sum_{k=1}^m r^{jk} x'_k$ , and for  $j = m + 1, \dots, n$ ,  $x_j = x'_j$ . With the above notations, we have an expression of the matrix  $\mathbf{A}_{ij}$  for  $1 \leq i \neq j \leq m$  in the following lemma.

**Lemma 2** For  $1 \leq i, j \leq m$ ,

$$\mathbf{A}_{ij} = a_i a_j (c_{ij} \Omega_0 + \sum_{l,k=1}^m c_{ij}^{lk} \mathbf{u}_l \mathbf{v}_k^T) \quad (33)$$

where

$$\Omega_0 = \sum_{k=m+1}^n \mathbf{u}_k \mathbf{u}_k^T = \mathbf{I} - \sum_{k=1}^m \mathbf{u}_k \mathbf{v}_k^T, \quad (34)$$

$$c_{ij} = E[\varphi'(w_i x'_i + b_i) \varphi'(w_j x'_j + b_j)], \quad (35)$$

$$c_{ij}^{lk} = E[\varphi'(w_i x'_i + b_i) \varphi'(w_j x'_j + b_j) (\sum_{s=1}^m r^{ls} x'_s) x'_k]. \quad (36)$$

The proof of Lemma 2 is given in Appendix 2.

**Remark 4** The vectors  $\{\mathbf{u}_{m+1}, \dots, \mathbf{u}_n\}$  are only used for theoretical analysis. They are not needed in the equation (33) because of the equation (34). To compute  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ , we define  $\mathbf{U}_1 = [\mathbf{u}_1, \dots, \mathbf{u}_m]$  and  $\mathbf{V}_1 = [\mathbf{v}_1, \dots, \mathbf{v}_m]$ . From  $\mathbf{U}_1^T \mathbf{V}_1 = \mathbf{I}$ , we then have

$$\mathbf{V}_1 = \mathbf{U}_1 (\mathbf{U}_1^T \mathbf{U}_1)^{-1} \quad (37)$$

which is the generalized inverse of  $\mathbf{U}_1$ .

**Remark 5** When  $i = j$ , the equation (33) gives another representation for  $\mathbf{A}_{ii}$  which is more complicated than the representation (23). But the new representation for  $\mathbf{A}_{ii}$  is needed to derive some properties of the Fisher information matrix. To compute the inverse of  $\mathbf{A}_{ii}$ , we use the equation (26). To compute the inverse of the Fisher information matrix, we also need to compute the inverse of matrices which have the same structure as  $\mathbf{A}_{ij}$  for  $i \neq j$ . This will be explained later in details when we give an algorithm to compute the inverse of the Fisher information matrix.

It follows from (32) that both  $x_i$  and  $x'_i$  are Gaussian random variables with a zero mean and

$$E[x_i x_j] = \mathbf{v}_i^T \mathbf{v}_j = r^{ij}, \quad E[x'_i x'_j] = \mathbf{u}_i^T \mathbf{u}_j = r_{ij}. \quad (38)$$

Therefore,

$$(x'_1, \dots, x'_n) \sim N(\mathbf{0}, \mathbf{R}), \quad (39)$$

$$(x_1, \dots, x_n) \sim N(\mathbf{0}, \mathbf{R}^{-1}), \quad (40)$$

and from (35) and (36)  $c_{ij}$  is a function of  $b_i, b_j, w_i, w_j$  and  $r_{ij}$ , and  $c_{ij}^{lk}$  is a function of  $b_i, b_j, w_i, w_j$  and  $(r_{ij})_{m \times m}$ .

**Lemma 3** For  $1 \leq i \leq m$ ,

$$\mathbf{A}_{i,m+1} = \mathbf{A}_{m+1,i}^T = \left( \sum_{k=1}^m c_{i1}^k \mathbf{v}_k, \dots, \sum_{k=1}^m c_{im}^k \mathbf{v}_k \right) \quad (41)$$

where

$$c_{ij}^k = E[\varphi'(w_i x'_i + b_i) \varphi(w_j x'_j + b_j) x'_k], \quad 1 \leq i, j, k \leq m. \quad (42)$$

$\mathbf{A}_{i,m+2}$  has the same structure as  $\mathbf{A}_{i,m+1}$ :

$$\mathbf{A}_{i,m+2} = \mathbf{A}_{m+2,i}^T = \left( \sum_{k=1}^m \tilde{c}_{i1}^k \mathbf{v}_k, \dots, \sum_{k=1}^m \tilde{c}_{im}^k \mathbf{v}_k \right) \quad (43)$$

where

$$\tilde{c}_{ij}^k = a_i a_j E[\varphi'(w_i x_i' + b_i) \varphi'(w_j x_j' + b_j) x_k'], \quad 1 \leq i, j, k \leq m.$$

$$\mathbf{A}_{m+1, m+1} = (b_{ij})_{m \times m} \quad (44)$$

with  $b_{ij} = E[\varphi(w_i x_i' + b_i) \varphi(w_j x_j' + b_j)]$  is a function of  $b_i, b_j, w_i, w_j$  and  $r_{ij}$ .

$$\mathbf{A}_{m+1, m+2} = \mathbf{A}_{m+2, m+1}^T = (\tilde{b}_{ij})_{m \times m} \quad (45)$$

with  $\tilde{b}_{ij} = a_j E[\varphi(w_i x_i' + b_i) \varphi'(w_j x_j' + b_j)]$ .

$$\mathbf{A}_{m+2, m+2} = (b'_{ij})_{m \times m} \quad (46)$$

with  $b'_{ij} = a_i a_j E[\varphi'(w_i x_i' + b_i) \varphi'(w_j x_j' + b_j)]$ .

The proof of Lemma 3 is given in Appendix 3.

**Remark 6** When the input is Gaussian random vector, the matrix  $\mathbf{G}$  can be calculated analytically by Lemma 1-3 except for some numerical integrations for those coefficients  $d_1, d_2$  and  $c_{ij}^{kl}$  etc. When the pdf of the input is unknown, we can estimate the Fisher information matrix based on the empirical distribution of the input. However, this is time consuming and also needs a large number of input examples. It is difficult to implement the natural gradient descent method as an on-line algorithm in this way. When the pdf of the input is known but different from the standard Gaussian distribution, either non-standard Gaussian or non-Gaussian, we need some preprocessing procedures.

### 3.4 Preprocessing

In the previous sections, the explicit form of  $\mathbf{G}$  is found by assuming a standard Gaussian input. This assumption facilitates the computation of  $\mathbf{G}$ . In fact, the explicit form of  $\mathbf{G}$  is still useful for non-Gaussian inputs if a preprocessing procedure is applied. This is equivalent to adding one preprocessing layer in front of the stochastic multi-layer perceptron defined by (1).

Assume the sampled input  $\mathbf{x}_t$  is an i.i.d. process. When the input is not a standard Gaussian process, we can use a linear or non-linear mapping to transform the input into a Gaussian process.

If the input is Gaussian but the covariance matrix is not an identity matrix, we can apply a linear transform  $\mathbf{u}_t = \mathbf{B} \mathbf{x}_t$  to obtain a standard Gaussian process  $\mathbf{u}_t$ . Here, the matrix  $\mathbf{B}$  can be found either by a batch algorithm or an on-line algorithm. By the batch algorithm, we first compute the sample covariance matrix

$$\hat{\mathbf{R}}_{\mathbf{x}} = \frac{1}{T-1} \sum_{t=1}^T (\mathbf{x}_t - \bar{\mathbf{x}})(\mathbf{x}_t - \bar{\mathbf{x}})^T,$$

$$\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t,$$

then compute the Cholesky factorization  $\hat{\mathbf{R}}_{\mathbf{x}} = \mathbf{F} \mathbf{F}^T$  and find  $\mathbf{B} = \mathbf{F}^{-1}$ . We can also apply the following on-line whitening algorithm in [6]:

$$\mathbf{B}_{t+1} = (1 + \mu) \mathbf{B}_t - \mu \mathbf{u}_t \mathbf{u}_t^T \mathbf{B}_t \quad (47)$$

where  $\mathbf{u}_t = \mathbf{B}_t \mathbf{x}_t$  and  $\mu > 0$  is a learning rate.

If the input is not Gaussian, we need a non-linear function to transform the input  $\mathbf{x}_t$  to a Gaussian process. Let  $F_{\mathbf{x}}(\cdot)$  be the cumulative distribution function (cdf) of  $\mathbf{x}_t$  and  $F_N(\cdot)$  be the cdf of n-dimensional Gaussian r.v., then we obtain a standard Gaussian process

$$\mathbf{u}_t = F_N^{-1}(F_{\mathbf{x}}(\mathbf{x}_t)) \sim N(\mathbf{0}, \mathbf{I}).$$

If  $F_{\mathbf{x}}(\cdot)$  is unknown, it can be approximated by the empirical distribution based on the data set. This method is applicable to the input with an arbitrary distribution but it is not applicable to on-line processing since the empirical distribution is computed by a batch algorithm. If the probability density function (pdf) of the input can be approximated by some expansions such as the Gram-Charlier expansion and the Edgeworth expansion [12],  $F_{\mathbf{x}}(\cdot)$  can be approximated on-line by using an adaptive algorithm to compute the moments or cumulants in these expansions.

After preprocessing, we use the data set  $\tilde{D}_T = \{(\mathbf{u}_t, z_t), t = 1, \dots, T\}$  instead of  $D_T$  to train the multi-layer perceptron.

## 4 Inverse of Fisher Information Matrix

Now we calculate the inverse of the Fisher information matrix which is needed in the natural gradient descent method.

We have already shown the explicit form (26) of the inverse of  $\mathbf{A}_{ii}$  for  $1 \leq i \leq m$ . To calculate the entire inverse of  $\mathbf{G}$ , we need to compute the inverse of a matrix with the similar representation as  $\mathbf{A}_{ij}$  for  $1 \leq i \neq j \leq m$ .

To this end, let us define

$$\begin{aligned} Gl(m, \mathfrak{R}) &= \{\mathbf{A} \in \mathfrak{R}^{m \times m} : \det(\mathbf{A}) \neq 0\}, \\ \mathcal{M} &= \{a_0 \Omega_0 + \sum_{i,j=1}^m a_{ij} \mathbf{u}_i \mathbf{v}_j^T : a_0 \neq 0, \mathbf{A} = (a_{ij}) \in Gl(m, \mathfrak{R})\}, \\ \overline{\mathcal{M}} &= \{a_0 \Omega_0 + \sum_{i,j=1}^m a_{ij} \mathbf{u}_i \mathbf{v}_j^T\} \quad (\text{the closure of } \mathcal{M}), \\ \widetilde{\mathcal{M}} &= \{[a_0, \mathbf{A}] : a_0 \neq 0, \mathbf{A} \in Gl(m, \mathfrak{R})\}, \end{aligned}$$

and recall the definition of  $\Omega_0 = \sum_{i,j=m+1}^n \mathbf{u}_i \mathbf{u}_j^T$ . The  $Gl(m, \mathfrak{R})$  is the set of all non-singular  $m \times m$  matrices. The set  $\widetilde{\mathcal{M}}$  is a product space of  $\mathfrak{R} \setminus \{0\}$  and  $Gl(m, \mathfrak{R})$ .

### 4.1 A group structure

Define a mapping  $\psi : \widetilde{\mathcal{M}} \rightarrow \mathfrak{R}^{n \times n}$  by

$$\psi(\tilde{A}) = a_0 \Omega_0 + \sum_{i,j=1}^m a_{ij} \mathbf{u}_i \mathbf{v}_j^T \in \mathcal{M}, \quad \text{for } \tilde{A} = [a_0, (a_{ij})] \in \widetilde{\mathcal{M}}.$$

Then  $\mathcal{M}$  is the image set of the mapping  $\psi$ .

For  $\tilde{A} = [a_0, \mathbf{A}]$  and  $\tilde{B} = [b_0, \mathbf{B}] \in \widetilde{\mathcal{M}}$ , we define a multiplication in  $\widetilde{\mathcal{M}}$  by

$$\tilde{A} \star \tilde{B} = [a_0 b_0, \mathbf{A} \mathbf{B}].$$

It is easy to verify that  $\widetilde{\mathcal{M}}$  is a group. For  $\tilde{A} = [a_0, \mathbf{A}] \in \widetilde{\mathcal{M}}$ , its inverse element is  $\tilde{A}^{-1} = [\frac{1}{a_0}, \mathbf{A}^{-1}]$ . The following lemma gives the relation between  $\mathcal{M}$  and  $\widetilde{\mathcal{M}}$ .

**Lemma 4** *The mapping  $\psi$  is an isomorphism between  $\mathcal{M}$  and  $\widetilde{\mathcal{M}}$  which implies that  $\mathcal{M}$  is sub-group of  $Gl(n, \mathbb{R})$  and for  $\mathbf{C} = a_0\Omega_0 + \sum_{i,j=1}^m a_{ij}\mathbf{u}_i\mathbf{v}_j^T \in \mathcal{M}$ ,*

$$\mathbf{C}^{-1} = \frac{1}{a_0}\Omega_0 + \sum_{i,j=1}^m a^{ij}\mathbf{u}_i\mathbf{v}_j^T \quad (48)$$

where the matrix  $(a^{ij})$  is the inverse of the matrix  $(a_{ij})^{-1}$ .

The proof of Lemma 4 is given in Appendix 4.

Regarding the relation between  $\mathcal{M}$  and  $\overline{\mathcal{M}}$ , we have

**Lemma 5**

$$\mathcal{M} = \overline{\mathcal{M}} \cap Gl(n, \mathbb{R}) \quad (49)$$

The proof of Lemma 5 is given in Appendix 5.

Let  $\mathbf{B}_{[m \times m]} = [\mathbf{B}_{ij}]_{m \times m}$  be a matrix consisting of  $m \times m$  blocks. To invert the Fisher information matrix, we need the following lemma regarding the inverse of a block matrix  $\mathbf{B}_{[m \times m]}$  whose blocks belong to  $\overline{\mathcal{M}}$ .

**Lemma 6** *If  $\mathbf{B}_{[m \times m]}$  is a positive definite matrix with the blocks  $\mathbf{B}_{ij} \in \overline{\mathcal{M}}$ , then the inverse  $\mathbf{B}_{[m \times m]}^{-1}$  has the same structure as  $\mathbf{B}_{[m \times m]}$ , namely,  $\mathbf{B}_{[m \times m]}^{-1} = [\mathbf{B}^{ij}]_{m \times m}$  in which  $\mathbf{B}^{ij} \in \overline{\mathcal{M}}$ .*

The proof of Lemma 6 is given here since it includes procedures to compute the inverse of a block matrix with blocks in  $\overline{\mathcal{M}}$ . These procedures will be referred to for several times in later sections.

Proof of Lemma 6:

We shall repeatedly apply the following 4-block matrix inverse formula in the current paper.

$$\begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{B}^{11} & \mathbf{B}^{12} \\ \mathbf{B}^{21} & \mathbf{B}^{22} \end{bmatrix}$$

where

$$\begin{aligned} \mathbf{B}^{11} &= \mathbf{B}_{11}^{-1} + \mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22,1}^{-1}\mathbf{B}_{21}\mathbf{B}_{11}^{-1}, \\ \mathbf{B}_{22,1} &= \mathbf{B}_{22} - \mathbf{B}_{21}\mathbf{B}_{11}^{-1}\mathbf{B}_{12}, \\ \mathbf{B}^{22} &= \mathbf{B}_{22,1}^{-1}, \text{ and} \\ \mathbf{B}^{12} &= (\mathbf{B}^{21})^T = -\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22,1}^{-1}. \end{aligned}$$

Since  $\mathbf{B}_{[m \times m]}$  is positive definite, the diagonal blocks  $\mathbf{B}_{ii}$  in  $\mathbf{B}_{[m \times m]}$  and  $\mathbf{B}_{[m \times m]}^{-1}$  are also positive definite. Let us first compute the inverse of  $[\mathbf{B}_{ij}]_{2 \times 2}$ . By Lemma 5, if an invertible matrix belongs to  $\overline{\mathcal{M}}$ , it also belongs to  $\mathcal{M}$ . So  $\mathbf{B}_{11} \in \mathcal{M}$ , and by Lemma 4  $\mathbf{B}_{11}^{-1} \in \mathcal{M}$ . Since the equality (65) in Appendix 4 still holds in  $\overline{\mathcal{M}}$ , the set  $\overline{\mathcal{M}}$  is closed under the matrix product. It is easy to see that  $\overline{\mathcal{M}}$  is also closed under the operations such as the matrix summation and scaling. By Lemma 5,  $\mathbf{B}_{22,1} \in \mathcal{M}$ , and so  $\mathbf{B}_{22,1}^{-1} \in \mathcal{M}$  again by Lemma 4. Hence,  $\mathbf{B}^{ij} \in \overline{\mathcal{M}}$  for  $i \neq j$ . The proposition of Lemma 6 is proved for  $m = 2$ . Assume the proposition is true for  $m = k - 1$ .

Again, we apply the 4-block matrix inverse formula to the following  $2 \times 2$  block matrix

$$\mathbf{B}_{[k \times k]} = \begin{bmatrix} \mathbf{B}'_{11} & \mathbf{B}'_{12} \\ \mathbf{B}'_{21} & \mathbf{B}'_{22} \end{bmatrix}$$

where

$$\begin{aligned} \mathbf{B}'_{11} &= \mathbf{B}_{[(k-1) \times (k-1)]} \\ \mathbf{B}'_{12} &= \begin{bmatrix} \mathbf{B}_{1,k} \\ \vdots \\ \mathbf{B}_{k-1,k} \end{bmatrix} \\ \mathbf{B}'_{21} &= [\mathbf{B}_{k,1}, \dots, \mathbf{B}_{k,k-1}] \\ \mathbf{B}'_{22} &= \mathbf{B}_{kk}. \end{aligned}$$

Because  $\overline{\mathcal{M}}$  is closed to the matrix product, summation and scaling,  $\mathbf{B}'_{22,1} = \mathbf{B}'_{22} - \mathbf{B}'_{21} \mathbf{B}'_{11}{}^{-1} \mathbf{B}'_{12}$  belongs to  $\overline{\mathcal{M}}$  by the inductive assumption that the proposition is true for  $m = k - 1$ . Since  $\mathbf{B}_{[k \times k]}$  is positive definite and so is  $\mathbf{B}'_{22,1}$  and  $\mathbf{B}'^{22} = \mathbf{B}'_{22,1}{}^{-1}$ . Similarly, we can prove that the blocks in  $\mathbf{B}'^{12}$ ,  $\mathbf{B}'^{21}$  and  $\mathbf{B}'^{11}$  belong to  $\overline{\mathcal{M}}$ . Therefore, we conclude that the proposition is also true for  $m = k$ .  
Q.E.D.

## 4.2 Inverse of $\mathbf{A}_{ij}$

Define matrices

$$\mathbf{C}_{ij} = \begin{bmatrix} c_{ij}^{1,1} & \cdots & c_{ij}^{1,m} \\ \cdots & \cdots & \cdots \\ c_{ij}^{m,1} & \cdots & c_{ij}^{m,m} \end{bmatrix}$$

for  $1 \leq i \neq j \leq m$ . Here,  $c_{ij}^{kl}$  are defined by (36). Applying Lemma 4-5, from Lemma 2 we have the following theorem.

**Theorem 1** For  $1 \leq i \neq j \leq m$ ,  $\mathbf{A}_{ij}$  is invertible iff  $c_{ij} \neq 0$  and the matrix  $\mathbf{C}_{ij}$  is invertible. When  $\mathbf{A}_{ij}$  is invertible, we have the following inverse formula

$$\mathbf{A}_{ij}^{-1} = \frac{1}{a_i a_j} \left( \frac{1}{c_{ij}} \Omega_0 + \sum_{l,k=1}^m d_{ij}^{lk} \mathbf{u}_l \mathbf{v}_k^T \right) \quad (50)$$

where  $(d_{ij}^{lk}) = \mathbf{C}_{ij}^{-1}$ .

## 4.3 Inverse of the whole Fisher information matrix

To compute the inverse of the Fisher information matrix  $\mathbf{G} = \frac{1}{\sigma^2} [\mathbf{A}_{ij}]_{(m+2) \times (m+2)}$ , we consider the following partition of  $\mathbf{G}$ :

$$\mathbf{G} = \begin{bmatrix} \mathbf{G}_{11} & \mathbf{G}_{12} \\ \mathbf{G}_{21} & \mathbf{G}_{22} \end{bmatrix} \quad (51)$$

where  $\mathbf{G}_{11}$  and  $\mathbf{G}_{22}$  are the  $m \times m$  and  $2 \times 2$  block matrices respectively on the left-upper and right-lower corners of  $\mathbf{G}$ . We first apply the 4-block matrix inverse formula to the above partition of  $\mathbf{G}$  and then apply the procedures in the proof of Lemma 6 to compute the inverse of  $\mathbf{G}_{11}$ .

#### 4.4 Complexity of computing the inverse of the Fisher information matrix

The method for computing  $\mathbf{G}^{-1}$  given in the above section has low complexity when  $m \ll n$ . We define the complexity of an algorithm by the time and space needed to run the algorithm. We count the time and space complexity by flops and units. A flop is one multiplication or one addition. A unit is the memory space for a variable.

One advantage of using the formulas (26) and (50) is to reduce the space complexity for storing  $\mathbf{G}$ . To store  $\mathbf{G}$  directly, we need  $\frac{1}{2}(mn + 2m)(mn + 2m + 1)$  units. But, by using the representation (33), when  $m \ll n$  the space complexity for storing  $\mathbf{G}$  is reduced to  $O(n)$  from  $O(n^2)$  since the dual bases are shared by all blocks in  $\mathbf{G}_{11}$ . Here,  $\mathbf{G}_{11}$  is defined in the partition (51) of  $\mathbf{G}$ .

The dual bases  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  are used in the theoretical analysis. The vectors  $\{\mathbf{u}_{m+1}, \dots, \mathbf{u}_n\}$  are not needed in the representation (33) of  $\mathbf{A}_{ij}$  and the inverse formula (50). Normalizing the weight vectors  $\{\mathbf{w}_1, \dots, \mathbf{w}_m\}$ , we obtain  $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$  with the time complexity  $O(n)$ . By the algorithm given in Remark 4, we can compute  $\mathbf{V}_1 = [\mathbf{v}_1, \dots, \mathbf{v}_m]$ . The time complexity is also  $O(n)$ .

In the partition (51), the dimension of  $\mathbf{G}_{11}$  is  $mn \times mn$  and that of  $\mathbf{G}_{22}$  is  $2m \times 2m$ . When  $m \ll n$ , the time complexity of computing  $\mathbf{G}^{-1}$  is determined by that of computing  $\mathbf{G}_{11}^{-1}$  which is computed by the procedures in the proof of Lemma 6 where the 4-block matrix inverse formula and Lemma 4-5 are repeatedly applied. In this process,  $\mathbf{G}_{11}^{-1}$  is recursively computed. The size of the working memory needed is still of order  $O(n^2)$  since  $m \ll n$ . But, when  $m \ll n$ , the time complexity of computing  $\mathbf{G}_{11}^{-1}$  is significantly reduced from  $O(n^3)$  by conventional methods such as the one in [11] to  $O(n^2)$  by our method. In summary, we have the following table to compare our method with the conventional methods:

Table 1: The complexity of our method and the conventional methods

Cost	Conventional methods	Our method
Space to store $\mathbf{G}$	$O(n^2)$	$O(n)$
Time to compute $\mathbf{G}^{-1}$	$O(n^3)$	$O(n^2)$
Working memory needed	$O(n^2)$	$O(n^2)$

Since the explicit expressions (26) and (50) are suitable for parallel matrix computations, we can further reduce the time complexity by using a computer with distributed memory multiprocessors.

#### 4.5 The Cramer-Rao lower bound for a committee machine

For a committee machine, we can further simplify the factor  $\text{Tr}(\mathbf{G}^{-1}(\boldsymbol{\theta}^*))$  in the Cramer-Rao lower bound inequality (3). When  $a_i = 1, b_i = 0, i = 1, \dots, m$ , the multi-layer perceptron (1) becomes a model for a committee machine. The dimension of the parameter space of this committee machine is  $nm$ .

The Fisher information matrix for the committee machine is

$$\mathbf{G}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} \mathbf{A}_{[m \times m]} = \frac{1}{\sigma^2} \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1m} \\ \cdots & \cdots & \cdots \\ \mathbf{A}_{m1} & \cdots & \mathbf{A}_{mm} \end{bmatrix}$$

where  $\boldsymbol{\theta} = (\mathbf{w}_1^T, \dots, \mathbf{w}_m^T)^T$  and  $\mathbf{A}_{ij}$  are the same as those for the multi-layer perceptron except  $a_i = 1$ .

By Lemma 6,

$$\mathbf{G}(\boldsymbol{\theta})^{-1} = \sigma^2 \mathbf{A}_{[m \times m]}^{-1} = \sigma^2 \begin{bmatrix} \mathbf{A}^{11} & \dots & \mathbf{A}^{1m} \\ \dots & \dots & \dots \\ \mathbf{A}^{m1} & \dots & \mathbf{A}^{mm} \end{bmatrix}$$

where  $\mathbf{A}^{ii} \in \mathcal{M}$  for all  $i$  and  $\mathbf{A}^{ij} \in \overline{\mathcal{M}}$  for  $i \neq j$ .

Note for  $\mathbf{A} = a_0 \Omega_0 + \sum_{k,l=1}^m a_{kl} \mathbf{u}_k \mathbf{v}_l^T \in \mathcal{M}$ ,

$$\text{Tr}(\mathbf{A}) = a_0(n - m) + \text{Tr}((a_{kl})) = a_0(n - m) + \sum_{k=1}^m a_{kk}$$

because of the properties (30) and (31).

Let  $\mathbf{A}^{ii} = a_0^i \Omega_0 + \sum_{k,l=1}^m a_{ii}^{kl} \mathbf{u}_k \mathbf{v}_l^T$ , then

$$\text{Tr}(\mathbf{A}^{ii}) = a_0^i(n - m) + \sum_{k=1}^m a_{ii}^{kk}.$$

Hence,

$$\text{Tr}(\mathbf{G}(\boldsymbol{\theta})^{-1}) = \sigma^2 \text{Tr}(\mathbf{A}_{[m \times m]}^{-1}) = \sigma^2 \sum_{i=1}^m \text{Tr}(\mathbf{A}^{ii}) = \sigma^2 (c_0(\boldsymbol{\theta}) + c_1(\boldsymbol{\theta})n)$$

where

$$\begin{aligned} c_1(\boldsymbol{\theta}) &= \sum_{i=1}^m a_0^i \\ c_0(\boldsymbol{\theta}) &= \sum_{i,k=1}^m a_{ii}^{kk} - c_1(\boldsymbol{\theta})m. \end{aligned}$$

So we have the following theorem which gives an order estimation of the Cramer-Rao lower bound for the committee machine.

**Theorem 2** *Assume  $m \ll n$  and  $\mathbf{G}(\boldsymbol{\theta})$  is non-singular around  $\boldsymbol{\theta}^*$ , then the Cramer-Rao lower bound inequality is*

$$E_{\boldsymbol{\theta}^*} [\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|^2] \geq \frac{\sigma^2}{t} (c_0 + c_1 n) \quad (52)$$

where  $c_0 = c_0(\boldsymbol{\theta}^*)$  and  $c_1 = c_1(\boldsymbol{\theta}^*)$  are two constants not depending on the input dimension  $n$ .

It is shown in [2] that the gradient descent learning rule (6) is Fisher efficient. In particular, for the committee machine, by Theorem 2 the natural gradient descent rule gives the mean square error

$$E_{\boldsymbol{\theta}^*} [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|^2] \approx \frac{\sigma^2}{t} (c_0 + c_1 n). \quad (53)$$

## 5 Natural gradient vector field vs. ordinary vector field

The behaviors of the two learning dynamics (5) and (6) are characterized, respectively, by the mean vector fields

$$\mathbf{v}_1(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}^*} \left[ \frac{\partial l}{\partial \boldsymbol{\theta}}(z|\mathbf{x}; \boldsymbol{\theta}) \right] = E_{\boldsymbol{\theta}^*} \left[ \begin{array}{c} \frac{\partial l}{\partial \mathbf{w}_1}(z|\mathbf{x}; \boldsymbol{\theta}) \\ \vdots \\ \frac{\partial l}{\partial \mathbf{w}_m}(z|\mathbf{x}; \boldsymbol{\theta}) \\ \frac{\partial l}{\partial \mathbf{a}}(z|\mathbf{x}; \boldsymbol{\theta}) \\ \frac{\partial l}{\partial \mathbf{b}}(z|\mathbf{x}; \boldsymbol{\theta}) \end{array} \right] \quad \text{and}$$

$$\mathbf{v}_2(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}^*} [\mathbf{G}^{-1}(\boldsymbol{\theta}) \frac{\partial l}{\partial \boldsymbol{\theta}}(z|\mathbf{x}; \boldsymbol{\theta})] = \mathbf{G}^{-1}(\boldsymbol{\theta}) \mathbf{v}_1(\boldsymbol{\theta}).$$

Again, the expectation with respect to  $p(\mathbf{x}, z; \boldsymbol{\theta}^*)$  is denoted by  $E_{\boldsymbol{\theta}^*}[\cdot]$  and the samples  $(\mathbf{x}, z)$  are independently drawn from the pdf  $p(\mathbf{x}, z; \boldsymbol{\theta}^*)$  for the teacher network:

$$z = \sum_{j=1}^m a_j^* \varphi(\mathbf{w}_j^{*T} \mathbf{x} + b_j^*) + \xi.$$

From (9)-(11), we have

$$E_{\boldsymbol{\theta}^*} \left[ \frac{\partial l}{\partial \mathbf{w}_i}(z|\mathbf{x}; \boldsymbol{\theta}) \right] = -\frac{1}{\sigma^2} E_{\boldsymbol{\theta}^*} \left[ \sum_{j=1}^m (a_j^* \varphi(\mathbf{w}_j^{*T} \mathbf{x} + b_j^*) - a_j \varphi(\mathbf{w}_j^T \mathbf{x} + b_j)) a_i \varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \mathbf{x} \right],$$

$$E_{\boldsymbol{\theta}^*} \left[ \frac{\partial l}{\partial \mathbf{a}}(z|\mathbf{x}; \boldsymbol{\theta}) \right] = -\frac{1}{\sigma^2} E_{\boldsymbol{\theta}^*} \left[ \sum_{j=1}^m (a_j^* \varphi(\mathbf{w}_j^{*T} \mathbf{x} + b_j^*) - a_j \varphi(\mathbf{w}_j^T \mathbf{x} + b_j)) \varphi(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \right],$$

$$E_{\boldsymbol{\theta}^*} \left[ \frac{\partial l}{\partial \mathbf{b}}(z|\mathbf{x}; \boldsymbol{\theta}) \right] = -\frac{1}{\sigma^2} E_{\boldsymbol{\theta}^*} \left[ \sum_{j=1}^m (a_j^* \varphi(\mathbf{w}_j^{*T} \mathbf{x} + b_j^*) - a_j \varphi(\mathbf{w}_j^T \mathbf{x} + b_j)) \{ \mathbf{a} \odot \varphi'(\mathbf{W}^T \mathbf{x} + \mathbf{b}) \} \right].$$

The above two equations are used to compute the vector field  $\mathbf{v}_1(\boldsymbol{\theta})$ . Note the  $\sigma^2$  in  $\mathbf{v}_1(\boldsymbol{\theta})$  and  $\mathbf{G}^{-1}(\boldsymbol{\theta})$  are canceled in computing  $\mathbf{v}_2(\boldsymbol{\theta})$ . So the vector field  $\mathbf{v}_2(\boldsymbol{\theta})$  does not depend on the noise level  $\sigma^2$ .

The following dynamic systems:

$$\frac{d\boldsymbol{\theta}}{dt} = -\eta \mathbf{v}_1(\boldsymbol{\theta}) \tag{54}$$

$$\frac{d\boldsymbol{\theta}}{dt} = -\eta \mathbf{v}_2(\boldsymbol{\theta}) = -\eta \mathbf{G}^{-1}(\boldsymbol{\theta}) \mathbf{v}_1(\boldsymbol{\theta}) \tag{55}$$

characterize the averaging behaviors of the gradient descent algorithm (5) and the natural gradient descent algorithm (6) respectively.

To show the difference between the two mean vector fields  $\mathbf{v}_1(\boldsymbol{\theta})$  and  $\mathbf{v}_2(\boldsymbol{\theta})$ , we consider a simple committee machine with 2-dimensional input, 2 hidden neurons,  $a_1 = a_2 = 1$  and the thresholds  $b_1 = b_2 = 0$ , and the 2 weight vectors  $\mathbf{w}_i$  with unit length  $\|\mathbf{w}_i\| = 1$ . The weight vectors are reparameterized by

$$\mathbf{w}_i = \begin{bmatrix} \cos(\alpha_i) \\ \sin(\alpha_i) \end{bmatrix}, \quad i = 1, 2.$$

Let  $\varphi(x) = \tanh(x)$ ,  $\sigma^2 = 1$  and  $\alpha_1^* = 0$  and  $\alpha_2^* = \frac{3\pi}{4}$  be the parameters for the teacher. Let  $\boldsymbol{\theta} = (\alpha_1, \alpha_2)$ . We calculate the mean vector fields  $v_1(\alpha_1, \alpha_2)$  and  $v_2(\alpha_1, \alpha_2)$  in the same way of deriving (54) and (55). The two mean vector fields are plotted in Figure 1-2.

Although the gradient descent system (54) and the natural gradient descent system (55) have the same set of fixed points, they generate different flows. The convergence of the ordinary gradient flow towards  $\boldsymbol{\theta}^*$  is slower than the natural gradient flow especially along the diagonal line in the  $(\alpha_1, \alpha_2)$ -space. On average, the learning dynamics of the natural gradient descent algorithm is optimized to achieve the Cramer-Rao lower bound.

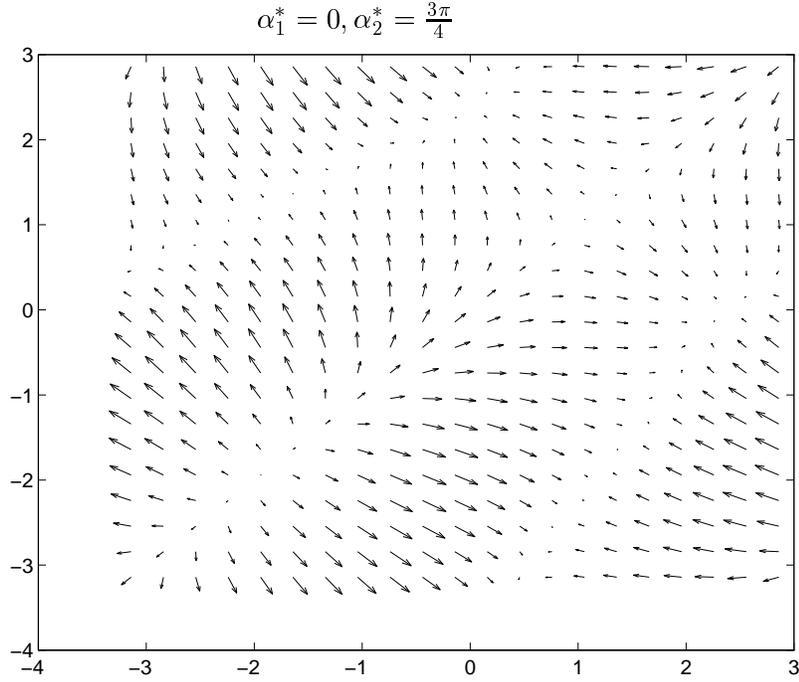


Figure 1: The mean vector field  $v_1(\alpha_1, \alpha_2)$

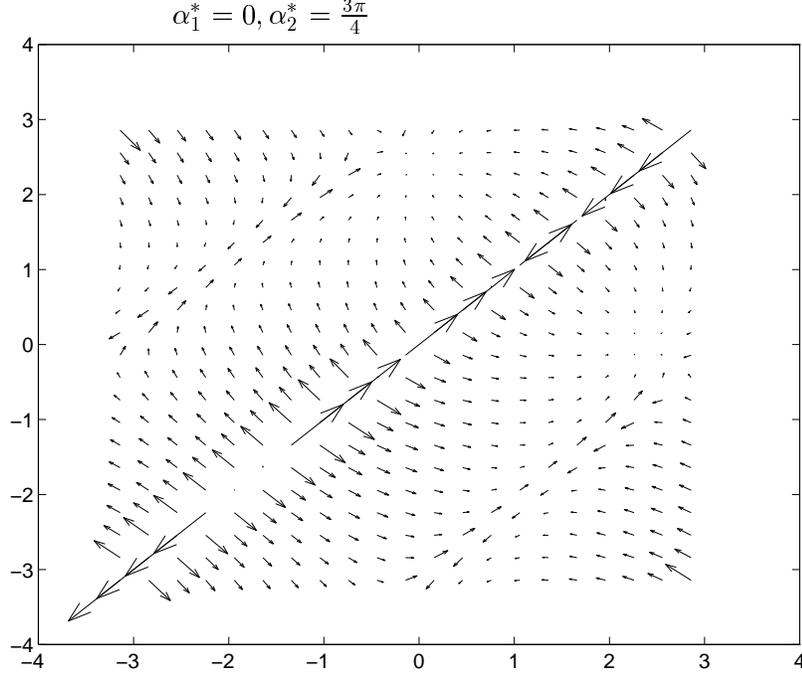


Figure 2: The mean vector field  $v_2(\alpha_1, \alpha_2)$

## 6 Simulation

### 6.1 Single-layer perceptron

Assume the input  $\mathbf{x}_t$  is of dimension  $n = 7$  subject to the Gaussian distribution  $\mathbf{x}_t \sim N(\mathbf{0}, \mathbf{I})$ . The output  $z_t$  are generated by the model:

$$z_t = \varphi(\mathbf{w}^{*T} \mathbf{x}_t) + \xi_t$$

where  $\varphi(u) = \frac{1-e^{-u}}{1+e^{-u}}$  and  $\xi_t \sim N(0, \sigma^2)$ . For the single-layer perceptron,  $z = \varphi(\mathbf{w}^T \mathbf{x})$ , we have the two on-line algorithms based on the gradient descent (GD) and the natural GD methods in the following forms:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu_0(t)(z_t - \varphi(\mathbf{w}_t^T \mathbf{x}_t))\varphi'(\mathbf{w}_t^T \mathbf{x}_t)\mathbf{x}_t \quad (56)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu_1(t)\mathbf{A}^{-1}(\mathbf{w}_t)(z_t - \varphi(\mathbf{w}_t^T \mathbf{x}_t))\varphi'(\mathbf{w}_t^T \mathbf{x}_t)\mathbf{x}_t \quad (57)$$

where

$$\mathbf{A}^{-1}(\mathbf{w}) = \frac{1}{d_1(\mathbf{w})}\mathbf{I} + \left(\frac{1}{d_2(\mathbf{w})} - \frac{1}{d_1(\mathbf{w})}\right)\frac{\mathbf{w}\mathbf{w}^T}{w^2}, \quad (58)$$

$$d_i(\mathbf{w}) = d_i(\mathbf{w}, 0), \quad i = 1, 2,$$

$d_i(\mathbf{w}, b)$  was defined previously by (24) and (25), and  $\mu_i(t)$  are two learning rate schedules defined by

$$\mu_i(t) = \eta_i \frac{1 + \frac{c_i}{\eta_i}t}{1 + \frac{c_i}{\eta_i}t + t^2}, \quad i = 0, 1. \quad (59)$$

Here  $t$  is the iteration index. The learning rate function  $\mu_i(t)$  is a special form of the following search-then-converge schedules proposed in [7]:

$$\mu(t) = \eta \frac{1 + \frac{c}{\eta} \frac{t}{\tau}}{1 + \frac{c}{\eta} \frac{t}{\tau} + \frac{t^2}{\tau}}. \quad (60)$$

$t < \tau$  is a “search phase” and  $t > \tau$  is a “converge phase”. The learning rate functions  $\mu_i(t)$  do not have the search phase but they start learning with a weaker converge phase when  $\eta_i$  are small. When  $t$  is large, each learning rate function  $\mu_i(t)$  decreases as  $\frac{c_i}{t}$ .

In this example, we choose  $\eta_0 = 1.25$ ,  $\eta_1 = 0.05$ ,  $c_0 = 8.75$ , and  $c_1 = 1$ . These parameters are selected by trial and error to optimize the performance of the GD and the natural GD methods at the noise level  $\sigma = 0.2$ . The learning rate functions  $\mu_0(t)$  and  $\mu_1(t)$  are compared in Figure 3 with their asymptotic curves  $\frac{c_i}{t}$ .

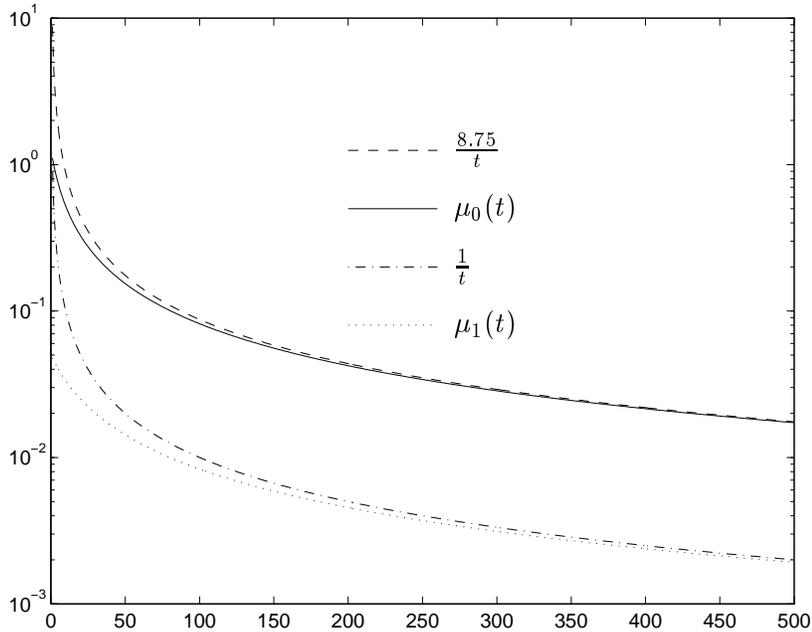


Figure 3:  $\mu_0(t)$  vs.  $\mu_1(t)$  when  $\eta_0 = 1.25$ ,  $\eta_1 = 0.05$ ,  $c_0 = 8.75$  and  $c_1 = 1$ .

A 7-dimensional vector is randomly chosen as  $\mathbf{w}^*$  for the teacher network. In this simulation, we choose

$$\mathbf{w}^* = [-1.1043, 0.4302, 1.1978, 1.5317, -2.2946, -0.7866, 0.4428]^T.$$

The training examples are the Gaussian random inputs to the teacher network and its outputs in Gaussian noise. Let  $\mathbf{w}_t$  and  $\tilde{\mathbf{w}}_t$  be the weight vectors driven by the equations (56) and (57) respectively. The error functions for the GD and the natural GD algorithms are  $\|\mathbf{w}_t - \mathbf{w}^*\|$  and  $\|\tilde{\mathbf{w}}_t - \mathbf{w}^*\|$ .

Let  $w^* = \|\mathbf{w}^*\|$ . From the equation (58), we obtain the Cramer-Rao Lower Bound (CRLB) for the deviation at the true weight vector  $\mathbf{w}^*$ :

$$\text{CRLB}(t) = \frac{\sigma}{\sqrt{t}} \sqrt{\frac{n-1}{d_1(w^*)} + \frac{1}{d_2(w^*)}}. \quad (61)$$

The error functions are compared with the CRLB in Figure 4.

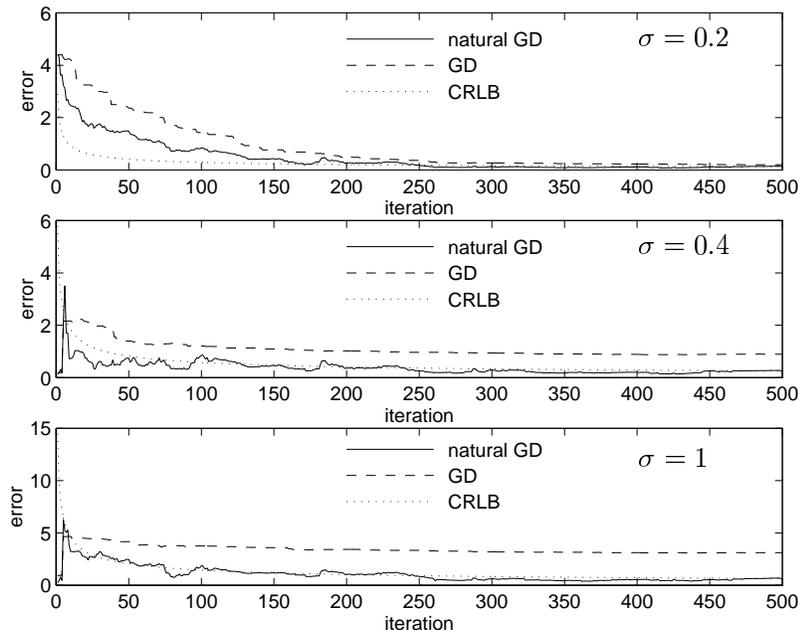


Figure 4: Performance of the GD and the natural GD at different noise levels  $\sigma = 0.2, 0.4, 1$ .

It is shown in Figure 4 that the natural GD algorithm reaches CRLB at different noise levels while the GD algorithm reaches the CRLB only at the noise level  $\sigma = 0.2$ . The robustness of the natural gradient descent against the additive noise in the training examples is clearly shown by this example.

## 6.2 Multi-layer perceptron

Let us reconsider the simple committee machine with 2-dimensional input and 2-hidden neurons given in Section 5. The problem is to train the following multi-layer perceptron:

$$y = \varphi(\mathbf{w}_1^T \mathbf{x}) + \varphi(\mathbf{w}_2^T \mathbf{x}) \quad (62)$$

based on the examples  $\{(\mathbf{x}_t, z_t), t = 1, \dots, T\}$  generated by the following stochastic committee machine:

$$z_t = \varphi(\mathbf{w}_1^{*T} \mathbf{x}_t) + \varphi(\mathbf{w}_2^{*T} \mathbf{x}_t) + \xi_t \quad (63)$$

with Gaussian noise  $\xi_t \sim N(0, \sigma^2)$  and Gaussian input  $\mathbf{x}_t \sim N(\mathbf{0}, \mathbf{I})$ . Assume  $\|\mathbf{w}_i^*\| = 1$ . We can reparameterize the weight vector to decrease the dimension of the parameter space from 4 to 2:

$$\mathbf{w}_i = \begin{bmatrix} \cos(\alpha_i) \\ \sin(\alpha_i) \end{bmatrix}, \quad \mathbf{w}_i^* = \begin{bmatrix} \cos(\alpha_i^*) \\ \sin(\alpha_i^*) \end{bmatrix}, \quad i = 1, 2.$$

The parameter space is  $\{\boldsymbol{\theta} = (\alpha_1, \alpha_2)\}$ . Assume the true parameters are  $\alpha_1^* = 0$  and  $\alpha_2^* = \frac{3\pi}{4}$ . Due to the symmetry, both  $\boldsymbol{\theta}_1^* = (0, \frac{3\pi}{4})$  and  $\boldsymbol{\theta}_2^* = (\frac{3\pi}{4}, 0)$  are true parameters. Let  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\theta}'_t$  be generated by the GD algorithm and the natural GD algorithm respectively. The errors are measured by

$$\begin{aligned}\varepsilon_t &= \min\{\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_1^*\|, \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_2^*\|\}, \\ \varepsilon'_t &= \min\{\|\boldsymbol{\theta}'_t - \boldsymbol{\theta}_1^*\|, \|\boldsymbol{\theta}'_t - \boldsymbol{\theta}_2^*\|\}.\end{aligned}$$

In this simulation, we first start the GD algorithm with the initial vector  $\boldsymbol{\theta}_0 = (0.1, 0.2)$  for 80 iterations, then we start the natural GD algorithm using the estimates of the parameters obtained by the GD algorithm at the 80-th iteration as the initial for the natural GD algorithm. The noise level is  $\sigma = 0.05$ .  $N$  independent runs are conducted to obtain  $\varepsilon_t(j)$  and  $\varepsilon'_t(j)$  for  $j = 1, \dots, N$ . Averaging these errors, we obtain

$$\bar{\varepsilon}_t = \sqrt{\frac{1}{N} \sum_{j=1}^N (\varepsilon_t(j))^2}, \quad \text{and} \quad \bar{\varepsilon}'_t = \sqrt{\frac{1}{N} \sum_{j=1}^N (\varepsilon'_t(j))^2}.$$

The averaging errors  $\bar{\varepsilon}_t$  and  $\bar{\varepsilon}'_t$  for the two algorithms are compared with the CRLB in Figure 5 for  $N = 10$ . The search-then-converge learning schedule (60) is used in the GD algorithm while the learning rate for the natural GD algorithm is simply the annealing rate  $\frac{1}{k}$ . The CRLB is given in Section 4.5 by Theorem 2.

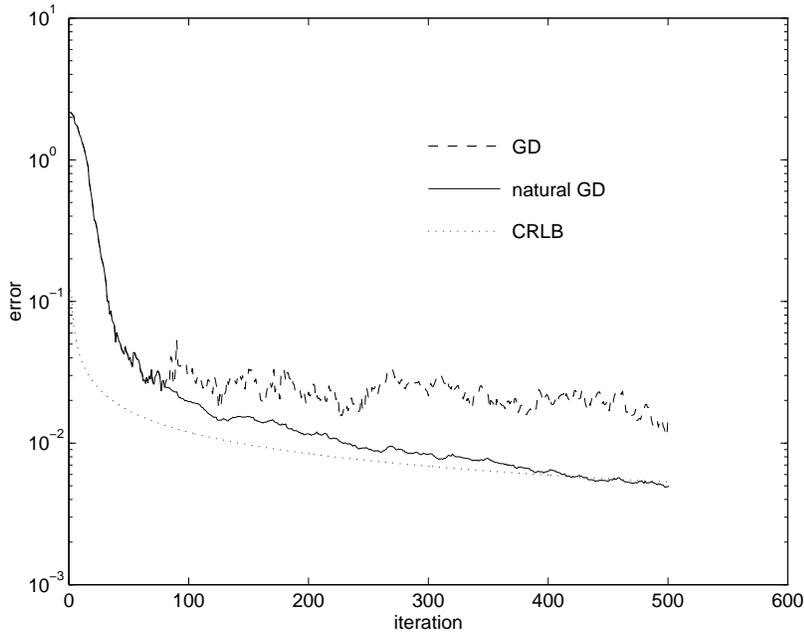


Figure 5: The GD vs. the natural GD

## 7 Conclusions

The natural gradient descent learning rule is statistically efficient. It can be used to train any adaptive system. But the complexity of this learning rule depends on the architecture of

the learning machine. The main difficulty in implementing this learning rule is to compute the inverse of the Fisher information matrix of a large size. For a multi-layer perceptron, we have shown an efficient scheme to represent the Fisher information matrix based on which the space for storing this large matrix is reduced to  $O(n)$  from  $O(n^2)$ . We have also shown an algorithm to compute the inverse of the Fisher information matrix. Its time complexity is of order  $O(n^2)$  when the input dimension  $n$  is much larger than the number of hidden neurons while the complexity of conventional algorithms for the same purpose is of order  $O(n^3)$ . An intermediate result is an explicit CRLB formula for a committee machine.

The simulation results have confirmed the fast convergence and statistical efficiency of the natural gradient descent learning rule. It is also verified that this learning rule is robust against the additive noise in the training examples.

## 8 Appendix

### Appendix 1 Proof of Lemma 1:

Let  $\mathbf{u}_1 = \mathbf{w}_i/w_i$  and extend it to  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  an orthogonal basis of  $\Re^n$ . The random input  $\mathbf{x}$  is decomposed as

$$\mathbf{x} = \sum_{j=1}^n c_j \mathbf{u}_j$$

where  $c_j = \mathbf{x}^T \mathbf{u}_j$ . Because of  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{I})$ ,  $c_j$  are i.i.d. with  $N(0, 1)$ . Noticing

$$\mathbf{x}\mathbf{x}^T = c_1^2 \mathbf{u}_1 \mathbf{u}_1^T + \sum_{j=2}^n c_j c_1 (\mathbf{u}_j \mathbf{u}_1^T + \mathbf{u}_1 \mathbf{u}_j^T) + \sum_{j,k=2}^n c_j c_k \mathbf{u}_j \mathbf{u}_k^T,$$

$E[c_j] = 0$  and  $E[c_j c_k] = \delta_{j,k}$ , we have

$$\begin{aligned} & E[(\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i))^2 \mathbf{x}\mathbf{x}^T] \\ &= E[(\varphi'(w_i c_1 + b_i))^2 c_1^2] \mathbf{u}_1 \mathbf{u}_1^T + E[(\varphi'(w_i c_1 + b_i))^2 c_1] \sum_{j=2}^n E[c_j] (\mathbf{u}_j \mathbf{u}_1^T + \mathbf{u}_1 \mathbf{u}_j^T) \\ &+ E[(\varphi'(w_i c_1 + b_i))^2] \sum_{j,k=2}^n E[c_j c_k] \mathbf{u}_j \mathbf{u}_k^T \\ &= E[(\varphi'(w_i c_1 + b_i))^2 c_1^2] \mathbf{u}_1 \mathbf{u}_1^T + E[(\varphi'(w_i c_1 + b_i))^2] \sum_{j=2}^n \mathbf{u}_j \mathbf{u}_j^T \\ &= d_2(w_i, b_i) \mathbf{u}_1 \mathbf{u}_1^T + d_1(w_i, b_i) \sum_{j=2}^n \mathbf{u}_j \mathbf{u}_j^T. \end{aligned}$$

Since  $\{\mathbf{u}_i\}$  are orthogonal,

$$\sum_{k=2}^n \mathbf{u}_k \mathbf{u}_k^T = \mathbf{I} - \mathbf{u}_1 \mathbf{u}_1^T.$$

Therefore,

$$E[(\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i))^2 \mathbf{x}\mathbf{x}^T] = d_1(w_i, b_i) \mathbf{I} + (d_2(w_i, b_i) - d_1(w_i, b_i)) \mathbf{u}_1 \mathbf{u}_1^T.$$

Q.E.D.

### Appendix 2 Proof of Lemma 2:

Here we use the notations introduced before Lemma 2.

Since  $\mathbf{u}_j = \mathbf{v}_j$  for  $j > m$ , we have  $x_j = x'_j$  for  $j > m$ , and  $E[x_i x'_j] = E[x_i x_j] = \delta_{i,j}$  for  $i, j > m$ .

When  $k \neq l$ ,  $E[x_k x'_l] = E[\mathbf{v}_k^T \mathbf{x} \mathbf{x}^T \mathbf{u}_l] = \mathbf{v}_k^T \mathbf{u}_l = 0$ . Since  $x_i$  and  $x'_i$  are Gaussian,  $x_k$  and  $x'_l$  are independent when  $k \neq l$ . Therefore,  $\{x'_1, \dots, x'_m\}$  and  $\{x_{m+1}, \dots, x_n\}$  are independent. Applying this result, we simplify the following

$$\begin{aligned}
& E[\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \varphi'(\mathbf{w}_j^T \mathbf{x} + b_j) \mathbf{x} \mathbf{x}^T] \\
&= E[\varphi'(w_i \mathbf{u}_i^T \mathbf{x} + b_i) \varphi'(w_j \mathbf{u}_j^T \mathbf{x} + b_j) (\sum_{l=1}^n x_l \mathbf{u}_l) (\sum_{k=1}^n x'_k \mathbf{v}_k^T)] \\
&= E[\varphi'(w_i x'_i + b_i) \varphi'(w_j x'_j + b_j) (\sum_{l,k=1}^m + \sum_{l=1}^m \sum_{k=m+1}^n + \sum_{l=m+1}^n \sum_{k=1}^m + \sum_{l,k=m+1}^n) x_l x'_k \mathbf{u}_l \mathbf{v}_k^T] \\
&= \sum_{l,k=1}^m c_{ij}^{lk} \mathbf{u}_l \mathbf{v}_k^T + c_{ij} \sum_{l,k=m+1}^n \delta_{l,k} \mathbf{u}_l \mathbf{v}_k^T + \text{Term1} + \text{Term2} \\
&= c_{ij} \sum_{k=m+1}^n \mathbf{u}_k \mathbf{u}_k^T + \sum_{l,k=1}^m c_{ij}^{lk} \mathbf{u}_l \mathbf{v}_k^T + \text{Term1} + \text{Term2}
\end{aligned}$$

where  $c_{ij}$  and  $c_{ij}^{lk}$  are defined by (35) and (36). Here,

$$\begin{aligned}
\text{Term1} &= \sum_{l=1}^m \sum_{k=m+1}^n E[\varphi'(w_i x'_i + b_i) \varphi'(w_j x'_j + b_j) x_l x_k] \mathbf{u}_l \mathbf{u}_k^T \quad \text{and} \\
\text{Term2} &= \sum_{l=m+1}^n \sum_{k=1}^m E[\varphi'(w_i x'_i + b_i) \varphi'(w_j x'_j + b_j) x_l x'_k] \mathbf{u}_l \mathbf{v}_k^T
\end{aligned}$$

are zero due to three reasons: 1)  $x_l = \sum_{s=1}^m r^{ls} x'_s$  for  $1 \leq l \leq m$ , 2)  $\{x'_1, \dots, x'_m\}$  and  $\{x_{m+1}, \dots, x_n\}$  are independent and 3)  $E[x_k] = 0$ . Finally, we obtain (33).

By the definition  $\mathbf{V} = (\mathbf{U}^{-1})^T$ , so  $\sum_{k=1}^n \mathbf{u}_k \mathbf{v}_k^T = \mathbf{U} \mathbf{V}^T = \mathbf{I}$  and

$$\Omega_0 = \sum_{k=m+1}^n \mathbf{u}_k \mathbf{u}_k^T = \mathbf{I} - \sum_{k=1}^m \mathbf{u}_k \mathbf{v}_k^T.$$

Q.E.D.

### Appendix 3 Proof of Lemma 3:

Again, we shall apply the properties that  $\{x'_1, \dots, x'_m\}$  and  $\{x_{m+1}, \dots, x_n\}$  are independent to simplify  $\mathbf{A}_{i,m+1}$  for  $1 \leq i \leq m$ :

$$\begin{aligned}
\mathbf{A}_{i,m+1} &= E[\varphi'(\mathbf{w}_i^T \mathbf{x} + b_i) \mathbf{x} \varphi(\mathbf{x}^T \mathbf{W} + \mathbf{b}^T)] \\
&= E[\varphi'(w_i x'_i + b_i) (\sum_{k=1}^m x'_k \mathbf{v}_k + \sum_{k=m+1}^n x'_k \mathbf{v}_k) (\varphi(w_1 x'_1 + b_1), \dots, \varphi(w_m x'_m + b_m))] \\
&= E[\varphi'(w_i x'_i) \sum_{k=1}^m x'_k \mathbf{v}_k (\varphi(w_1 x'_1 + b_1), \dots, \varphi(w_m x'_m + b_m))] \\
&\quad + E[\varphi'(w_i x'_i + b_i) \sum_{k=m+1}^n x'_k \mathbf{v}_k (\varphi(w_1 x'_1 + b_1), \dots, \varphi(w_m x'_m + b_m))]. \tag{64}
\end{aligned}$$

The first term on the right hand side (RHS) of the above equation (64) becomes

$$(\sum_{i=1}^m c_{i1}^k \mathbf{v}_k, \dots, \sum_{i=1}^m c_{im}^k \mathbf{v}_k)$$

where  $c_{ij}^k$  are defined by (42). Since  $E[x_k] = 0$ , the second term on the RHS of (64) is zero. Denote

$$\tilde{\varphi}(\mathbf{u}^T) = \mathbf{a}^T \odot \varphi'(\mathbf{u}^T) = [a_1, \dots, a_m \varphi'(u_m)].$$

Similar to the derivation of  $\mathbf{A}_{i,m+1}$ , we obtain the expression (64) for  $\mathbf{A}_{i,m+2}$ .

It is easy to obtain  $\mathbf{A}_{m+1,m+1} = (b_{ij})_{m \times m}$  where

$$b_{ij} = E[\varphi(\mathbf{w}_i^T \mathbf{x} + b_i) \varphi(\mathbf{w}_j^T \mathbf{x} + b_j)] = E[\varphi(w_i x'_i + b_i) \varphi(w_j x'_j + b_j)].$$

$b_{ij}$  are functions of  $b_i$ ,  $b_j$ ,  $w_i$ ,  $w_j$ , and  $r_{ij}$  since  $E[x'_i x'_j] = r_{ij}$  for  $1 \leq i, j \leq m$ . Similarly, we obtain (45) and (46).

Q.E.D.

**Appendix 4** *Proof of Lemma 4:*

For  $\tilde{A} = [a_0, (a_{ij})] \in \tilde{\mathcal{M}}$ , define

$$\psi(\tilde{A}) = a_0\Omega_0 + \sum_{i,j=1}^m a_{ij}\mathbf{u}_i\mathbf{v}_j^T.$$

To show that  $\psi : \tilde{\mathcal{M}} \rightarrow \mathcal{M}$  is one-to-one, we assume  $\tilde{B} = [b_0, (b_{ij})] \in \tilde{\mathcal{M}}$  and  $\psi(\tilde{A}) = \psi(\tilde{B})$ , i.e.,

$$a_0\Omega_0 + \sum_{i,j=1}^m a_{ij}\mathbf{u}_i\mathbf{v}_j^T = b_0\Omega_0 + \sum_{i,j=1}^m b_{ij}\mathbf{u}_i\mathbf{v}_j^T.$$

Since  $\Omega_0\mathbf{u}_i\mathbf{v}_j^T = 0$  for  $1 \leq i, j \leq m$ , from the above equation we first obtain  $a_0 = b_0$  then  $a_{ij} = b_{ij}$  after applying the property (31). Therefore,  $\psi$  is a one-to-one mapping from  $\tilde{\mathcal{M}}$  to  $\mathcal{M}$ . Let  $\mathbf{C} = \psi(\tilde{A})$  and  $\mathbf{C}' = \psi(\tilde{B})$ , then

$$\begin{aligned} \mathbf{C}\mathbf{C}' &= a_0b_0\Omega_0 + \sum_{i,j,l,k=1}^m a_{ij}b_{lk}\mathbf{u}_i\mathbf{v}_j^T\mathbf{u}_l\mathbf{v}_k^T \\ &= a_0b_0\Omega_0 + \sum_{i,j,l,k=1}^m a_{ij}b_{lk}\delta_{j,l}\mathbf{u}_i\mathbf{v}_k^T = a_0b_0\Omega_0 + \sum_{i,j,k=1}^m a_{ij}b_{jk}\mathbf{u}_i\mathbf{v}_k^T \\ &= a_0b_0\Omega_0 + \sum_{i,k=1}^m d_{ik}\mathbf{u}_i\mathbf{v}_k^T, \end{aligned} \tag{65}$$

where  $(d_{ij}) = (a_{ij})(b_{ij})$ . So  $\mathbf{C}\mathbf{C}' = \psi([a_0b_0, (a_{ij})(b_{ij})])$ , i.e.,

$$\psi(\tilde{A})\psi(\tilde{B}) = \psi(\tilde{A} \star \tilde{B}).$$

Hence,  $\psi$  is an isomorphic mapping between  $\tilde{\mathcal{M}}$  and  $\mathcal{M}$ , and  $\mathcal{M}$  is a sub-group of  $Gl(n, \mathbb{R})$ . For any element in  $\mathcal{M}$  with the form  $\mathbf{C} = a_0\Omega_0 + \sum_{i,j=1}^m a_{ij}\mathbf{u}_i\mathbf{v}_j^T = \psi([a_0, (a_{ij})])$ , to obtain the inverse of  $\mathbf{C}$ , we first compute the inverse in  $\tilde{\mathcal{M}}$

$$\tilde{A}^{-1} = \left[ \frac{1}{a_0}, (a_{ij})^{-1} \right]$$

then the inverse in  $\mathcal{M}$

$$\mathbf{C}^{-1} = \psi(\tilde{A}^{-1}) = \frac{1}{a_0}\Omega_0 + \sum_{i,j=1}^m a^{ij}\mathbf{u}_i\mathbf{v}_j^T$$

by the relation  $(\psi(\tilde{A}))^{-1} = \psi(\tilde{A}^{-1})$ , where  $(a^{ij}) = (a_{ij})^{-1}$ .  
Q.E.D.

**Appendix 5** *Proof of Lemma 5:*

It follows from Lemma 4:

$$\mathcal{M} \subseteq \overline{\mathcal{M}} \cap Gl(n, \mathbb{R}).$$

Let  $\mathbf{A} \in \overline{\mathcal{M}}$  and  $\det(\mathbf{A}) \neq 0$ . From  $\mathbf{A} = a_0\Omega_0 + \sum_{i,j=1}^m a_{ij}\mathbf{u}_i\mathbf{v}_j^T$ , we have

$$\begin{aligned} \mathbf{A}\mathbf{u}_k &= a_0\mathbf{u}_k \text{ for } k = m+1, \dots, n, \text{ and} \\ \mathbf{A}\mathbf{u}_k &= [\mathbf{u}_1, \dots, \mathbf{u}_m][a_{1k}, \dots, a_{mk}]^T \text{ for } k = 1, \dots, m. \end{aligned}$$

So  $\mathbf{A}[\mathbf{u}_{m+1}, \dots, \mathbf{u}_n] = a_0[\mathbf{u}_{m+1}, \dots, \mathbf{u}_n]$  and  $\mathbf{A}[\mathbf{u}_1, \dots, \mathbf{u}_m] = [\mathbf{u}_1, \dots, \mathbf{u}_m]\mathbf{B}$  where  $\mathbf{B} = (a_{ij})$ . Hence,

$$\mathbf{AU} = \mathbf{A}[\mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{u}_{m+1}, \dots, \mathbf{u}_n] = \mathbf{U} \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & a_0\mathbf{I} \end{bmatrix}.$$

Therefore,  $\det(\mathbf{B})a_0^{n-m} = \det(\mathbf{AU}) \neq 0$  which implies  $\det(\mathbf{B}) \neq 0$  and  $a_0 \neq 0$ . So  $\mathbf{A} \in \mathcal{M}$ . Q.E.D.

## References

- [1] S. Amari. *Differential-Geometrical Methods in Statistics, Lecture Notes in Statistics vol.28*. Springer, 1985.
- [2] S. Amari. Natural gradient works efficiently in learning. *Accepted by Neural Computation*, 1997.
- [3] S. Amari. Neural learning in structured parameter spaces – natural Riemannian gradient. In *Advances in Neural Information Processing Systems, 9, MIT Press: Cambridge, MA. (to appear)*, 1997.
- [4] S. Amari and N. Murata. Statistical theory of learning curves under entropic loss criterion. *Neural Computation*, 5:140–153, 1993.
- [5] S. Amari, N. Murata, K.-R. Müller, M. Finke, and H. H. Yang. Asymptotic statistical theory of overtraining and cross-validation. *IEEE Trans. on Neural Networks*, page to appear, 1997.
- [6] J.-F. Cardoso and B. Laheld. Equivariant adaptive source separation. *IEEE Trans. on Signal Processing*, 44(12):3017–3030, December 1996.
- [7] C. Darken and J. Moody. Towards faster stochastic gradient search. In *Advances in Neural Information Processing Systems, 4, eds. Moody, Hanson, and Lippmann, Morgan Kaufmann, San Mateo*, pages 1009–1016, 1992.
- [8] R. A. Jacobs. Increased rates of convergence through learning rate adaptation. *Neural Networks*, 1:295–307, 1988.
- [9] M. K. Murray and J. W. Rice. *Differential Geometry and Statistics*. New York: Chapman & Hall, 1993.
- [10] B. D. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.
- [11] G. W. Stewart. *Introduction to Matrix Computations*. New York: Academic Press, 1973.
- [12] A. Stuart and J. K. Ord. *Kendall's Advanced Theory of Statistics*. Edward Arnold, 1994.