# Complexity Issues in Natural Gradient Descent Method for Training Multilayer Perceptrons

**Howard Hua Yang**
*Department of Computer Science, Oregon Graduate Institute, Portland, OR 97291, U.S.A.*

**Shun-ichi Amari**
*Laboratory for Information Synthesis, RIKEN Brain Science Institute, Saitama 351-01, Japan*

The natural gradient descent method is applied to train an $n$-$m$-1 multilayer perceptron. Based on an efficient scheme to represent the Fisher information matrix for an $n$-$m$-1 stochastic multilayer perceptron, a new algorithm is proposed to calculate the natural gradient without inverting the Fisher information matrix explicitly. When the input dimension $n$ is much larger than the number of hidden neurons $m$, the time complexity of computing the natural gradient is $O(n)$.

## 1 Introduction

Amari (1997, 1998) has shown that the natural gradient descent learning rule is statistically efficient. In principle, this learning rule can be used to train any adaptive system, but the complexity of this rule depends on the architecture of the learning machine. The main difficulty in implementing the rule is designing a fast algorithm to compute the natural gradient. For an $n$-$m$-1 multilayer perceptron, we show in this article that the natural gradient can be computed in $O(n)$ flops when $m \ll n$. Here, a flop is a floating-point operation: an add or a multiplication.

Orr and Leen (1997) used the curvature information (the Hessian matrix) in the nonlinear adaptive momentum scheme to optimize the convergence of the stochastic gradient descent. They have shown that the complexity of their algorithm is $O(n)$. However, their algorithm is different from the natural gradient descent. Instead of using the momentum scheme, we use the inverse of the Fisher information matrix to transform the stochastic gradient to optimize the learning dynamics. This method is based on statistical inference. A review of the learning of artificial neural networks is given by Yang, Murata, and Amari (1998).

The rest of this article is organized as follows. The stochastic multilayer perceptron is described in section 2. The natural gradient descent method using the inverse of the Fisher information matrix and its properties are

discussed in section 3. A constructive procedure to compute the inverse of the Fisher information matrix is formulated in section 4. For a single-layer perceptron, we calculate an explicit expression of the natural gradient in section 4.1. It is obvious from this expression that the time complexity for computing the natural gradient for the single-layer perceptron is $O(n)$. Importantly, this is also true for a multilayer perceptron, but its proof is not straightforward. To prove this result, we analyze the structure of the Fisher information matrix in section 4.2. Based on this analysis, we discuss the time complexities of computing the inverse of the Fisher information matrix and the natural gradient in section 5 for the multilayer perceptron. We first discuss the simplest multilayer perceptron, a committee machine, in section 5.1, and then discuss the multilayer perceptron in section 5.2. The conclusions are summarized in section 6.

## 2 Stochastic Multilayer Perceptron

We consider the following stochastic multilayer perceptron model:

$$z = \sum_{i=1}^{m} a_i \varphi(w_i^T x + b_i) + \xi, \tag{2.1}$$

where $(\cdot)^T$ denotes the transpose, $\xi \sim N(0, \sigma^2)$ is the gaussian additive noise with the variance $\sigma^2$, and $\varphi(x)$ is a differentiable output function of hidden neurons. Assume that the multilayer network has an $n$-dimensional input, $m$ hidden neurons, a one-dimensional output, and $m \leq n$. Denote $a = (a_1, \cdots, a_m)^T$ the weight vector of the output neuron, $w_i = (w_{1i}, \cdots, w_{ni})^T$ the weight vector of the $i$th hidden neuron, and $b = (b_1, \cdots, b_m)^T$ the vector of thresholds for the hidden neurons. Let $W = [w_1, \cdots, w_m]$ be a matrix formed by column weight vectors $w_i$; then equation 2.1 can be rewritten as

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

Here, the scalar function $\varphi$ operates on each component of the vector $W^T x + b$.

## 3 Fisher Information Matrix and the Natural Gradient Descent

The stochastic perceptron model enables us to derive learning algorithms for neural networks from some statistical inference methods and evaluate the performance of the algorithms systematically.

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

$$p(x, z; W, a, b) = p(z|x; W, a, b)p(x),$$

where $p(z|x; W, a, b)$ is the conditional pdf of the output $z$ given the input

$x$ and $p(x)$ is the pdf of the input $x$. The loss function is defined as the negative log-likelihood function

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

where

$$\theta = (w_1^T, \cdots, w_m^T, a^T, b^T)^T, \quad \text{and}$$

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

Since $p(x)$ does not depend on $\theta$, minimizing the loss function $L(x, z; \theta)$ is equivalent to minimizing $l(z|x; \theta)$.

Given the training set

$$D_T = \{(x_t, z_t), \ t = 1, \cdots, T\},$$

minimizing the loss function $L(D_T; \theta) = \prod_{t=1}^{T} L(x_t, z_t; \theta)$ is equivalent to minimizing the training error

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

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

$$G = G(\theta) = E_\theta\left[\frac{\partial L}{\partial \theta}(\frac{\partial L}{\partial \theta})^T\right] = E_\theta\left[\frac{\partial l}{\partial \theta}(\frac{\partial l}{\partial \theta})^T\right], \tag{3.1}$$

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

Let $\Theta$ be a parameter set and $\mathcal{P} = \{p(x, z; \theta) : \ \theta \in \Theta\}$ be a family of pdfs parameterized by $\theta \in \Theta$. The Kullback-Leibler divergence between two pdfs in $\mathcal{P}$ is

$$D(\theta, \theta') = \int dx dz \, p(x, z; \theta) \log \frac{p(x, z; \theta)}{p(x, z; \theta')}.$$

Let $G(\theta)$ be the Riemannian metric tensor for the Riemannian space $\Theta$. The squared length of a small $d\theta$ in $\Theta$ is

$$\|d\theta\|^2_{G(\theta)} = d\theta^T G(\theta)d\theta.$$

The two spaces $\mathcal{P}$ and $\Theta$ are naturally related by the following equation:

$$D(\theta, \theta') = \frac{1}{2}d\theta^T G(\theta)d\theta + O(\|d\theta\|^3),$$

where $\theta' = \theta + d\theta$ and $d\theta$ is small.

Amari (1998) showed that in the Riemannian space $\Theta$, for any differentiable loss function $F(\theta)$, the steepest descent direction is given by

$$-G^{-1}(\theta)\frac{\partial F}{\partial \theta}.$$

Based on this result, Amari (1997) proposed the following natural gradient descent algorithm:

$$\theta_{t+1} = \theta_t - \frac{\mu}{t}G^{-1}(\theta)\frac{\partial F}{\partial \theta},$$

where $\mu$ is a learning rate. In particular, when negative log-likelihood function is chosen as the loss function, the above algorithm is known as the method of scoring in statistics. Amari (1998) showed that this algorithm gives a Fisher efficient online estimator; the asymptotic variance of $\theta_t$ driven by equation 3.5 satisfies

$$E[(\theta_t - \theta^*)(\theta_t - \theta^*)^T \mid \theta^*] \approx \frac{1}{t}G^{-1}(\theta^*), \tag{3.2}$$

which gives the mean square error

$$E[\|\theta_t - \theta^*\|^2 \mid \theta^*] \approx \frac{1}{t}\mathrm{Tr}(G^{-1}(\theta^*)). \tag{3.3}$$

The above property is verified by the simulation results in Yang and Amari (1998).

Yang and Amari (1997a) have shown that $G(\theta) = \frac{1}{\sigma^2}A(\theta)$ where $A(\theta)$ does not depend on $\sigma^2$ (the variance of the additive noise).

Define

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

and

$$\widetilde{\nabla}l_1(z|x; \theta) = A^{-1}(\theta)\frac{\partial l_1}{\partial \theta}(z|x; \theta). \tag{3.4}$$

To minimize the training error, we propose the following natural gradient descent algorithm:

$$\theta_{t+1} = \theta_t - \frac{\mu}{t}\widetilde{\nabla}l_1(z_t|x_t; \theta_t), \tag{3.5}$$

where the natural gradient $\widetilde{\nabla}l_1(z_t|x_t; \theta_t)$ does not depend on $\sigma^2$.

To implement this algorithm, we need to compute the natural gradient $\widetilde{\nabla}l_1(z_t|x_t; \theta_t)$ in each iteration. If we compute the inverse of $A(\theta_t)$ first, then the natural gradient, the time complexity is $O(n^3)$.

We proposed a method (Yang & Amari, 1997b) based on the conjugate gradient method to compute the natural gradient without inverting the Fisher information matrix. The idea is to solve the following linear equation,

$$A(\theta_t)y = \frac{\partial l_1}{\partial \theta}(z_t|x_t; \theta_t),$$

for $y = \widetilde{\nabla} l_1(z_t|x_t; \theta_t)$ without inverting the matrix $A(\theta_t)$. Since the matrix $A(\theta)$ is $(n+2)m \times (n+2)m$, it will take at most $(n+2)m$ steps by the conjugate gradient method to compute $\widetilde{\nabla} l_1(z_t|x_t; \theta_t)$ and each step needs $O(n)$ flops to compute a matrix vector product. The number of flops needed to compute the natural gradient by the conjugate gradient method is $O(n^2)$ when $m \ll n$ and $O(n^3)$ when $m = O(n)$. Using the conjugate gradient method to compute the natural gradient is also useful for other probability families when the method of scoring is used to find the maximum likelihood estimates (MLEs). However, the conjugate gradient method is still slow in computing the natural gradient for training the stochastic perceptron. We shall explore the structure of the Fisher information matrix of the stochastic perceptron and develop an algorithm that computes the natural gradient in just $O(n)$ flops.

## 4 Computing the Natural Gradient for a Stochastic Perceptron

To find an analytic form of the Fisher information matrix, we assume a white gaussian input $x \sim N(0, I)$ with an identity matrix $I$ as the covariance matrix of the input. Yang and Amari (1997a) proposed an efficient scheme to represent the Fisher information matrix for the multilayer perceptron. Due to this scheme, when $m \ll n$, the storage space needed for the Fisher information matrix is $O(n)$ units rather than $O(n^2)$ units, where a unit is the memory space to keep a parameter or a variable. Based on this scheme, we found a constructive procedure to compute the inverse of the Fisher information matrix with $O(n^2)$ flops. This procedure can be improved to compute the natural gradient with only $O(n)$ flops.

**4.1 Single-Layer Perceptron.** We shall give some explicit expressions of the Fisher information matrix and the natural gradient for a one-layer stochastic perceptron.

Since $z = \varphi(w^T x + b) + \xi$, we have

$$\frac{\partial l}{\partial w} = -\frac{\xi}{\sigma^2}\varphi'(w^T x + b)x,$$

$$\frac{\partial l}{\partial b} = -\frac{\xi}{\sigma^2}\varphi'(w^T x + b).$$

The Fisher information matrix is

$$G(\theta) = \frac{1}{\sigma^2}A(\theta) = \frac{1}{\sigma^2}\begin{bmatrix} A_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

where $\theta = (w^T, b)^T$,

$$A_{11} = E[(\varphi'(w^T x + b))^2 x x^T],$$
$$a_{12} = a_{21}^T = E[(\varphi'(w^T x + b))^2 x],$$
$$a_{22} = E[(\varphi'(w^T x + b))^2].$$

Let $w = \|w\|$ and $u_1 = w/w$ and extend $u_1$ to an orthogonal basis $\{u_1, \cdots, u_n\}$ for $\Re^n$. The gaussian input $x$ is decomposed as $x = \sum_i x_i u_i$; we have $x_i \sim N(0, 1)$ and

$$E[x_i x_j] = E[u_i^T x x^T u_j] = u_i^T u_j = \delta_{ij}, \tag{4.1}$$

since $x \sim N(0, I)$. The outer product $x x^T$ can be written as

$$x x^T = x_1^2 u_1 u_1^T + \sum_{j=2}^{n} x_j x_1 (u_j u_1^T + u_1 u_j^T) + \sum_{j,k=2}^{n} x_j x_k u_j u_k^T. \tag{4.2}$$

Noticing $w^T x = w x_1$ and applying equations 4.1 and 4.2 we have,

$$A_{11} = E[(\varphi'(w^T x + b))^2 x x^T]$$
$$= E[(\varphi'(w x_1 + b))^2 x_1^2] u_1 u_1^T$$
$$\quad + E[(\varphi'(w x_1 + b))^2 x_1] \sum_{j=2}^{n} E[x_j](u_j u_1^T + u_1 u_j^T)$$
$$\quad + E[(\varphi'(w x_1 + b))^2] \sum_{j,k=2}^{n} E[x_j x_k] u_j u_k^T$$
$$= E[(\varphi'(w x_1 + b))^2 x_1^2] u_1 u_1^T + E[(\varphi'(w x_1 + b))^2] \sum_{j=2}^{n} u_j u_j^T.$$

Let us define three integrals:

$$d_0(w, b) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\varphi'(w x + b))^2 e^{-\frac{x^2}{2}} dx > 0, \tag{4.3}$$

$$d_1(w, b) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\varphi'(w x + b))^2 x e^{-\frac{x^2}{2}} dx, \tag{4.4}$$

$$d_2(w, b) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\varphi'(w x + b))^2 x^2 e^{-\frac{x^2}{2}} dx > 0. \tag{4.5}$$

Then we have

$$A_{11} = d_2(w, b) u_1 u_1^T + d_0(w, b) \sum_{j=2}^{n} u_j u_j^T,$$

$$a_{12} = E[(\varphi'(w x_1 + b))^2 x_1] u_1 = d_1(w, b) u_1, \quad \text{and}$$

$$a_{22} = d_0(w, b).$$

$A_{11}$ can be written as

$$A_{11} = d_0(w, b)I + (d_2(w, b) - d_0(w, b))u_1 u_1^T,$$

since $\{u_i\}$ are orthogonal and

$$\sum_{k=2}^{n} u_i u_i^T = I - u_1 u_1^T.$$

Summarizing previous calculations for the blocks in $G(\theta)$, we have

$$A_{11} = d_0(w, b)I + (d_2(w, b) - d_0(w, b))ww^T/w^2, \tag{4.6}$$

$$a_{12} = a_{21}^T = d_1(w, b)w/w, \quad a_{22} = d_0(w, b). \tag{4.7}$$

Saad and Solla (1995) used the scaled error function

$$\varphi(x) = \text{erf}(\frac{x}{\sqrt{2}}) = \frac{2}{\sqrt{\pi}} \int_0^{\frac{x}{\sqrt{2}}} e^{-t^2} dt$$

as the sigmoid function for the hidden units in order to obtain an analytic expression of the generalization error. If we choose $\varphi(x)$ as the sigmoid function for the hidden neurons, we obtain the following close forms for the integrals 4.3 through 4.5:

$$d_0(w, b) = \frac{\sqrt{2}}{\pi \sqrt{w^2 + 0.5}} \exp\{-b^2 + \frac{b^2 w^2}{w^2 + 0.5}\},$$

$$d_2(w, b) = d_0(w, b)\frac{1}{w^2 + 0.5}(\frac{1}{2} + \frac{w^2 b^2}{w^2 + 0.5}),$$

$$d_1(w, b) = -\frac{\sqrt{2}wb}{\pi (w^2 + 0.5)^{3/2}} \exp\{-b^2 + \frac{b^2 w^2}{w^2 + 0.5}\}.$$

Even for the single-layer perceptron, the size of the Fisher information matrix is $(n + 1) \times (n + 1)$. However, the Fisher information matrix can be generated online with $O(n^2)$ flops by equations 4.6 and 4.7. If one wants to trace the Fisher information matrix, one need only $O(n)$ units to store $w$ and $b$ in each iteration.

To compute the inverse of $A(\theta)$, we need the following well-known inverse formula of a four-block matrix:

**Lemma 1.**
$$\begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}^{-1} = \begin{bmatrix} B^{11} & B^{12} \\ B^{21} & B^{22} \end{bmatrix}$$

*provided* $|B_{11}| \neq 0$ *and* $|B_{22} - B_{21}B_{11}^{-1}B_{12}| \neq 0$. *Here,*
$B^{11} = B_{11}^{-1} + B_{11}^{-1}B_{12}B_{22,1}^{-1}B_{21}B_{11}^{-1},$
$B_{22,1} = B_{22} - B_{21}B_{11}^{-1}B_{12},$
$B^{22} = B_{22,1}^{-1},$
*and* $B^{12} = (B^{21})^T = -B_{11}^{-1}B_{12}B_{22,1}^{-1}.$

Applying the above inverse formula of a four-block matrix, we obtain

$$A^{-1}(\boldsymbol{\theta}) = \begin{bmatrix} \boldsymbol{A}^{11} & \boldsymbol{a}^{12} \\ \boldsymbol{a}^{21} & a^{22} \end{bmatrix}$$

where

$$\boldsymbol{A}^{11} = \frac{1}{d_0}\boldsymbol{I} + \left(\frac{1}{d_2} - \frac{1}{d_0}\right)\frac{\boldsymbol{w}\boldsymbol{w}^T}{w^2} + \frac{d_1^2}{d_2(d_0d_2 - d_1^2)}\frac{\boldsymbol{w}\boldsymbol{w}^T}{w^2}, \tag{4.8}$$

$$\boldsymbol{a}^{12} = (\boldsymbol{a}^{21})^T = -\frac{d_1d_2}{d_2(d_0d_2 - d_1^2)}\boldsymbol{w}, \tag{4.9}$$

$$a^{22} = \frac{d_2}{d_0d_2 - d_1^2}. \tag{4.10}$$

By using equations 4.8 and 4.9, the time complexity for computing $A^{-1}(\boldsymbol{\theta})$ is $O(n^2)$. If we compute the inverse of the Fisher information matrix and the function $l_1(z|\boldsymbol{x}; \boldsymbol{\theta})$ in two separate steps and then multiply them together, the time complexity for computing the natural gradient is $O(n^2)$. Instead of computing the natural gradient in this way, we compute it in a single step by applying equations 4.8 and 4.9:

$$\widetilde{\nabla}l_1(z|\boldsymbol{x}; \boldsymbol{\theta}) = -h(\boldsymbol{x}, z; \boldsymbol{w}, b) \tag{4.11}$$

$$\times \begin{bmatrix} \frac{\boldsymbol{x}}{d_0} + (\frac{1}{d_2} - \frac{1}{d_0} + \frac{d_1^2}{d_2(d_0d_2 - d_1^2)})\frac{\boldsymbol{w}^T\boldsymbol{x}}{w^2}\boldsymbol{w} - \frac{d_1d_2}{d_2(d_0d_2 - d_1^2)}\boldsymbol{w} \\ -\frac{d_1d_2\boldsymbol{w}^T\boldsymbol{x}}{d_2(d_0d_2 - d_1^2)} + \frac{d_2}{d_0d_2 - d_1^2} \end{bmatrix},$$

where $h(\boldsymbol{x}, z; \boldsymbol{w}, b) = (z - \varphi(\boldsymbol{w}^T\boldsymbol{x} + b))\varphi'(\boldsymbol{w}^T\boldsymbol{x} + b)$. It needs only $O(n)$ flops to compute the natural gradient $\widetilde{\nabla}l_1(z|\boldsymbol{x}; \boldsymbol{\theta})$ by equation 4.11.

**4.2 Multilayer Perceptron.** The main difficulty in implementing the natural gradient descent algorithm (see equation 3.5) is to compute the natural gradient online. The number of parameters in the stochastic perceptron defined by equation 2.1 is $m(n+2)$, so the size of the Fisher information matrix is $m(n + 2) \times m(n + 2)$. In this article, we assume $m \ll n$. To compute the inverse of an $n \times n$ matrix, it generally needs $O(n^3)$ flops by commonly used methods such as the lower and upper triangular factorization together with the triangular matrix inverse (Stewart, 1973) .

In the previous section, we showed some procedures for computing the Fisher information matrix and the natural gradient for a one-layer perceptron. Following similar procedures, Yang and Amari (1997a) found a scheme to represent the Fisher information matrix for multilayer perceptrons. Based on this scheme, we gave a method that requires $O(n^2)$ flops to compute the inverse of the Fisher information matrix. In this article, we adapt this method

to compute the natural gradient $\widetilde{\nabla} l_1(z|x; \theta)$. We shall show that the adapted method requires only $O(n)$ flops to compute the natural gradient.

We first briefly describe this scheme to represent the Fisher information matrix and then introduce the method to compute the natural gradient.

We need some notations for block matrix. An $m \times n$ block matrix $X$ is denoted by $X = [X_{ij}]_{[m \times n]}$ or $X = [X_{ij}]_{i=1,\cdots,m, j=1,\cdots,n}$.

Let

$$A(\theta) = [A_{ij}]_{[(m+2) \times (m+2)]}$$

be a partition of $A(\theta)$ corresponding to the partition of

$$\theta = (w_1^T, \cdots, w_m^T, a^T, b^T)^T.$$

Define

$$u_i = w_i / \|w_i\|, i = 1, \cdots, m,$$
$$U_1 = [u_1, \cdots, u_m], \text{ and}$$
$$[v_1, \cdots, v_m] = U_1(U_1^T U_1)^{-1}.$$

Let $r_{ij} = u_i^T u_j$, $R_1 = (r_{ij})_{m \times m}$ and $R_1^{-1} = (r^{ij})_{m \times m}$. Define $x_i' = u_i^T x$ for $i = 1, \cdots, m$; then

$$(x_1', \cdots, x_m') \sim N(0, R_1).$$

The structure of the blocks in $A(\theta)$ is elucidated by the following two lemmas proved in Yang and Amari (1997a), which gives an efficient scheme to represent the Fisher information matrix for the stochastic multilayer perceptron.

**Lemma 2.** *For $1 \le i, j \le m$, and $a_i \ne 0$,*

$$A_{ij} = a_i a_j (c_{ij} \Omega_0 + \sum_{l,k=1}^{m} c_{ij}^{lk} u_l v_k^T) \tag{4.12}$$

*where*

$$\Omega_0 = \sum_{k=m+1}^{n} u_k u_k^T = I - \sum_{k=1}^{m} u_k v_k^T, \tag{4.13}$$

$$c_{ij} = E[\varphi'(w_i x_i' + b_i)\varphi'(w_j x_j' + b_j)], \tag{4.14}$$

$$c_{ij}^{lk} = E\left[ \varphi'(w_i x_i' + b_i)\varphi'(w_j x_j' + b_j) \left( \sum_{s=1}^{m} r^{ls} x_s' \right) x_k' \right]. \tag{4.15}$$

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

$$A_{i,m+1} = A_{m+1,i}^T = \left( \sum_{k=1}^{m} c_{i1}^k v_k, \cdots, \sum_{k=1}^{m} c_{im}^k v_k \right), \tag{4.16}$$

*where*

$$c_{ij}^k = E[\varphi'(w_i x_i' + b_i)\varphi(w_j x_j' + b_j)x_k'], \quad 1 \leq i, j, k \leq m. \tag{4.17}$$

*$A_{i,m+2}$ has the same structure as $A_{i,m+1}$:*

$$A_{i,m+2} = A_{m+2,i}^T = \left( \sum_{k=1}^{m} \tilde{c}_{i1}^k v_k, \cdots, \sum_{k=1}^{m} \tilde{c}_{im}^k v_k \right), \tag{4.18}$$

*where*

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

$$A_{m+1,m+1} = (b_{ij})_{m \times m} \tag{4.19}$$

*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 $r_{ij}$.*

$$A_{m+1,m+2} = A_{m+2,m+1}^T = (\tilde{b}_{ij})_{m \times m}, \tag{4.20}$$

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

$$A_{m+2,m+2} = (b_{ij}')_{m \times m} \tag{4.21}$$

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

Assume $\varphi(x) = \text{erf}(x/\sqrt{2})$ where $\text{erf}(u)$ is the error function. It is not difficult to obtain analytic expressions for $c_{ij}, b_{ij}', c_{ij}^{lk}$, and $\tilde{c}_{ij}^k$. When $b_i = 0, i = 1, \cdots, m$, applying the multivariate gaussian average method in Saad and Solla (1995), we can obtain analytic expressions for these coefficients: $c_{ij}^k, b_{ij}$, and $\tilde{b}_{ij}$. But when $b_i \neq 0$, the analytic expressions for these coefficients seemingly do not exist. These coefficients are $m$-dimensional integrals, which can be approximated offline by some numerical integral methods.

Recalling that $n$ is the input dimension and $m$ is the number of hidden neurons, in this article, we assume $m \ll n$. Using lemmas 2 and 3, instead of using $O(n^2)$ units to store the Fisher information matrix, we only need $O(n)$ units to store the vectors $u_k$ and $v_k$. $O(n^2)$ flops are needed to construct the Fisher information matrix from these vectors.

## 5 Complexity Issues

We shall discuss the time complexity to compute $A(\theta)$, $A^{-1}(\theta)$ and the natural gradient for the multilayer perceptron. By lemmas 2 and 3, to compute the blocks in $A(\theta)$, we need to compute $u_k$ and $v_k$, $k = 1, \cdots, m$. It needs $O(n)$ flops to compute $u_k$ and $v_k$ by their definitions. The time complexity to compute $A(\theta)$ is $O(n^2)$ since each outer product $u_k v_k^T$ needs $O(n^2)$ flops. It is shown in Yang and Amari (1997a) that about $O(n^2)$ flops are needed to compute the inverse of the Fisher information matrix. If we compute the Fisher information matrix first and then compute the natural gradient, this will take $O(n^2)$ flops in each iteration to implement the natural gradient descent. This is in the same order as the conjugate gradient method proposed by Yang and Amari (1997b). In the rest of this article, we show that the time complexity of computing the natural gradient for the multilayer perceptron can be significantly improved by using the representation scheme in lemmas 2 and 3 and by computing the natural gradient without inverting the Fisher information matrix. We shall show that we can compute the natural gradient directly with only $O(n)$ flops. To show this, we reexamine the process of computing the inverse of the Fisher information matrix.

We need the following notations:

$$Gl(m, \Re) = \{A \in \Re^{m \times m} : \det(A) \neq 0\}.$$

$$\mathcal{M} = \left\{ a_0 \Omega_0 + \sum_{i,j=1}^{m} a_{ij} u_i v_j^T : a_0 \neq 0, A = (a_{ij}) \in Gl(m, \Re) \right\}.$$

$$\overline{\mathcal{M}} = \left\{ a_0 \Omega_0 + \sum_{i,j=1}^{m} a_{ij} u_i v_j^T \right\} \quad \text{(the closure of } \mathcal{M}\text{)}.$$

$$\widetilde{\mathcal{M}} = \{[a_0, A] : a_0 \neq 0, A \in Gl(m, \Re)\}.$$

$$\overline{\widetilde{\mathcal{M}}} = \{[a_0, A] : a_0 \in \Re, A \in \Re^{m \times m}\} \text{( the closure of } \widetilde{\mathcal{M}} \text{ )}.$$

$\widetilde{\mathcal{M}}$ is a subset of $\overline{\widetilde{\mathcal{M}}}$. An element in $\overline{\widetilde{\mathcal{M}}}$ is denoted by $[a_0, A]$ or $[a_0, (a_{ij})]$. Define a mapping $\psi : \overline{\widetilde{\mathcal{M}}} \to \Re^{n \times n}$,

$$\psi(\widetilde{A}) = a_0 \Omega_0 + \sum_{i,j=1}^{m} a_{ij} u_i v_j^T \in \overline{\mathcal{M}}, \quad \text{for } \widetilde{A} = [a_0, (a_{ij})] \in \overline{\widetilde{\mathcal{M}}}.$$

For $\widetilde{A} = [a_0, A]$ and $\widetilde{B} = [b_0, B] \in \overline{\widetilde{\mathcal{M}}}$, we define multiplication, summation, and scaling in $\overline{\widetilde{\mathcal{M}}}$ by

$$\widetilde{A} \star \widetilde{B} = [a_0 b_0, AB],$$
$$\widetilde{A} + \widetilde{B} = [a_0 + b_0, A + B],$$
$$a[a_0, A] = [aa_0, aA].$$

It is easy to verify that $\widetilde{\mathcal{M}}$ is a group. For $\widetilde{A} = [a_0, A] \in \widetilde{\mathcal{M}}$, its inverse element is $\widetilde{A}^{-1} = [\frac{1}{a_0}, A^{-1}]$.

Note that $\mathcal{M}$ is the image set of $\widetilde{\mathcal{M}}$ under the mapping $\psi$. It is not difficult to show that $\psi : \widetilde{\mathcal{M}} \to \mathcal{M}$ is an isomorphism, meaning that $\psi$ is one-to-one and $\psi(\widetilde{A})\psi(\widetilde{B}) = \psi(\widetilde{A} \star \widetilde{B})$ for $\widetilde{A}$ and $\widetilde{B} \in \widetilde{\mathcal{M}}$. This implies that $\mathcal{M}$ is a subgroup of $Gl(n, \Re)$ and for $C = a_0 \Omega_0 + \sum_{i,j=1}^{m} a_{ij} u_i v_j^T \in \mathcal{M}$,

$$C^{-1} = \frac{1}{a_0} \Omega_0 + \sum_{i,j=1}^{m} a^{ij} u_i v_j^T \tag{5.1}$$

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

The following lemma, proved in Yang and Amari (1997a), gives the relations between $\mathcal{M}$ and $\overline{\mathcal{M}}$.

**Lemma 4.**

$$\mathcal{M} = \overline{\mathcal{M}} \bigcap Gl(n, \Re). \tag{5.2}$$

We shall apply equation 5.1 and lemma 4 to study the structure of $A^{-1}(\theta)$.

**5.1 Committee Machine.** Let us first consider a committee machine that is a multilayer perceptron with $a_i = 1$ and $b_i = 0$ for $i = 1, \cdots, m$.

Repeatedly applying the inverse formula of the four-block matrix, formula 5.1, and lemma 4, we know that $A^{-1}(\theta)$ has the same structure as $A(\theta)$: $A^{-1}(\theta) = [A^{ij}]_{i=1,\cdots,m, j=1,\cdots,m}$, where every $A^{ij} \in \overline{\mathcal{M}}$. In this process, all matrix multiplications, summations, and scalings are carried out in the space $\overline{\widetilde{\mathcal{M}}}$, and all matrix inversions are made in the group $\widetilde{\mathcal{M}}$. Therefore, the time cost for computing the representations of $A^{ij}$ in $\overline{\widetilde{\mathcal{M}}}$ is $O(m^4)$, which does not depend on $n$. Only when each block $A^{ij}$ is realized in the matrix space $\Re^{n \times n}$, the time complexity of computing $A^{-1}(\theta)$ is $O(n^2)$ since realizing each $u_k v_k^T$ needs $n^2$ flops of computations. We shall avoid computing the outer product $u_k v_k^T$ in computing the natural gradient $\widetilde{\nabla} l_1$.

Let

$$\widetilde{\nabla} l_1 = \begin{bmatrix} (\widetilde{\nabla} l_1)_{w_1} \\ \vdots \\ (\widetilde{\nabla} l_1)_{w_m} \end{bmatrix}$$

be a partition of $\widetilde{\nabla} l_1$ corresponding to the partition of $\boldsymbol{\theta} = (w_1^T, \cdots, w_m^T)^T$. Here each $(\widetilde{\nabla} l_1)_{w_i}$ is an $n \times 1$ vector.

From the block representation of $A^{-1}(\boldsymbol{\theta})$, we have

$$(\widetilde{\nabla} l_1)_{w_i} = \sum_{j=1}^{m} A^{ij} \frac{\partial l_1}{\partial w_j}.$$

Noticing $A^{ij} \in \overline{\mathcal{M}}$, we have

$$A^{ij} = b_{ij} \Omega_0 + \sum_{k,l=1}^{m} b_{ij}^{kl} u_k v_l^T,$$

where the coefficients $b_{ij}$ and $b_{ij}^{kl}$ are determined in the process of computing $A^{-1}(\boldsymbol{\theta})$. From the above expression of $A^{ij}$, we obtain

$$(\widetilde{\nabla} l_1)_{w_i} = \sum_{j=1}^{m} \left[ b_{ij} \frac{\partial l_1}{\partial w_j} - b_{ij} \sum_{k=1}^{m} \left( v_k^T \frac{\partial l_1}{\partial w_j} \right) u_k + \sum_{k,l=1}^{m} b_{ij}^{kl} \left( v_l^T \frac{\partial l_1}{\partial w_j} \right) u_k \right]. \quad (5.3)$$

It requires only the vector inner product, the vector addition, and no matrix multiplication to compute the above expression, so the time complexity for computing the natural gradient $\widetilde{\nabla} l_1$ is $O(n)$.

In summary, for the committee machine, $O(n^2)$ flops due to the vector outer product and the matrix addition are required to compute the inverse of the Fisher information matrix, and $O(n)$ flops due to vector inner product and vector addition are needed to compute the natural gradient by formula 5.3. In the next section, we show that this conclusion is also true for the multilayer perceptron.

### 5.2 Multilayer Perceptron. Recalling that

$$A(\boldsymbol{\theta}) = [A_{ij}]_{[(m+2) \times (m+2)]}$$

is a partition of $A(\boldsymbol{\theta})$ corresponding to the partition of

$$\boldsymbol{\theta} = (w_1^T, \cdots, w_m^T, a^T, b^T)^T,$$

we consider a new partition of $A(\boldsymbol{\theta})$ as:

$$A(\boldsymbol{\theta}) = \begin{bmatrix} B_{11} & B_{12} \\ B_{12}^T & B_{22} \end{bmatrix},$$

where

$$B_{11} = [A_{ij}]_{i=1,\cdots,m, j=1,\cdots,m},$$

$$B_{12} = [A_{ij}]_{i=1,\cdots,m, j=m+1,m+2},$$

$$B_{22} = [A_{ij}]_{i=m+1,m+2, j=m+1,m+2}.$$

Assuming that both $B_{11}$ and $B_{22} - B_{21}B_{11}^{-1}B_{12}$ are nonsingular and applying the four-block inverse formula, we have

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

where $B^{ij}$ are defined in the same way as lemma 1.

To characterize the structure of these blocks $B^{ij}$ in $A^{-1}(\theta)$ we need some notations and lemmas.

Define

$$\mathcal{M}_u = \{A \in \Re^{n \times m} : \text{ each column of } A \text{ belongs to } R(U_1)\},$$

$$\mathcal{M}_v = \{A \in \Re^{n \times m} : \text{ each column of } A \text{ belongs to } R(V_1)\}.$$

Here, $R(U_1)$ is the range of $U_1$ generated by all linear combinations of the columns in $U_1$. $R(V_1)$ is defined in the same way from $V_1$.

**Lemma 5.** *If $A \in \overline{\mathcal{M}}$ and $B \in \mathcal{M}_v$, then $AB \in \mathcal{M}_u$ and $AB$ has the following form,*

$$\left[ \sum_{i=1}^{m} a_1^i u_i, \cdots, \sum_{i=1}^{m} a_m^i u_i \right],$$

*where the coefficients $a_j^i$ are computed in $O(n)$ flops.*

The proof of lemma 5 is given in appendix A.

**Lemma 6.** *Let $A \in \mathcal{M}_u$ and $B \in \mathcal{M}_v$; then the complexity for computing the $m \times m$ matrix $B^T A$ is $O(n)$.*

The proof of lemma 6 is given in appendix B.

**Lemma 7.** *If $A \in \mathcal{M}_{\mathcal{U}}$ and $B \in \Re^{m \times m}$, then $AB \in \mathcal{M}_{\mathcal{U}}$. $AB$ has the following form,*

$$\left[ \sum_{i=1}^{m} b_1^i u_i, \cdots, \sum_{i=1}^{m} b_m^i u_i \right],$$

*where the coefficients $b_j^i$ are computed in $m^3$ flops.*

The proof of lemma 7 is given in appendix C.

Let us recall the definition of $B^{11}$:

$$B^{11} = B_{11}^{-1} + B_{11}^{-1} B_{12} B_{22,1}^{-1} B_{21} B_{11}^{-1}, \tag{5.4}$$
$$B_{22,1} = B_{22} - B_{21} B_{11}^{-1} B_{12}.$$

The matrix $B_{11}$ has the same structure as the Fisher information matrix for the committee machine, so $B_{11}^{-1}$ is an $m \times m$ matrix with each of its blocks belonging to $\overline{\mathcal{M}}$. Therefore, the time complexity for computing $B_{11}^{-1}$ is the same as that for the committee machine. However, we shall show below that the blocks in $B_{11}^{-1} B_{12} B_{22,1}^{-1} B_{21} B_{11}^{-1}$ do not belong to $\overline{\mathcal{M}}$, and the time complexity for computing these blocks needs to be discussed separately.

To discuss the time complexity of computing the blocks $B^{ij}$, we need the following lemma:

**Lemma 8.** *The $m \times 2$ block matrix $B_{11}^{-1} B_{12}$ has the form*

$$B_{11}^{-1} B_{12} = [X_{ij}]_{[m \times 2]},$$

*where*

$$X_{ij} = \left[ \sum_{k=1}^{m} x_{ij}^{k1} u_k, \cdots, \sum_{k=1}^{m} x_{ij}^{km} u_k \right], \tag{5.5}$$

*and $O(n)$ flops are needed to compute the coefficients $x_{ij}^{kl}$.*

The proof of lemma 8 is given in appendix D.

**Lemma 9.** *It takes $O(n)$ flops to compute the $2m \times 2m$ matrix $B_{22,1}$.*
$B_{11}^{-1} B_{12} B_{22,1}^{-1} B_{21} B_{11}^{-1} = [Z_{ij}]_{[m \times m]}$ *is an $m \times m$ block matrix and each block $Z_{ij}$*

*has the following form:*

$$Z_{ij} = \sum_{k,l=1}^{m} z_{ij}^{kl} u_k u_l^T.$$

*It takes $O(n)$ flops to compute these coefficients $z_{ij}^{kl}$ in the above expression.*

The proof of lemma 9 is given in appendix E.

**Theorem 1.**   *The complexity for computing $B^{22} = B_{22,1}^{-1}$ is $O(n)$. The $m \times 2$ block matrix $B^{12}$ has the following form:*

$$B^{12} = -B_{11}^{-1} B_{12} B_{22,1}^{-1} = -[Y_{ij}]_{[m \times 2]},$$

*where*

$$Y_{ij} = \left[ \sum_{k=1}^{m} y_{ij}^{k1} u_k, \cdots, \sum_{k=1}^{m} y_{ij}^{km} u_k \right], \tag{5.6}$$

*and $O(n)$ flops are needed to compute the coefficients $y_{ij}^{kl}$.*
   *The $m \times m$ block matrix $B^{11}$ has the following form:*

$$B^{11} = [\widetilde{B}_{ij}]_{[m \times m]},$$

*where*

$$\widetilde{B}_{ij} = a_{ij} \Omega_0 + \sum_{k,l=1}^{m} b_{kl}^{ij} u_k v_l^T + \sum_{k,l=1}^{m} z_{kl}^{ij} u_k u_l^T, \tag{5.7}$$

*and $O(n)$ flops are needed to compute these coefficients $a_{ij}$, $b_{kl}^{ij}$ and $z_{kl}^{ij}$.*

**Proof.** By (1) in lemma 9, it takes $O(n)$ flops to compute $B_{22,1}$, which is of $2m \times 2m$. Therefore, $O(n)$ flops are needed to compute $B_{22,1}^{-1}$. This is shown in the second part of the proof for lemma 9. By equation 5.4 and the second statement in lemma 9, $B^{11}$ is an $m \times m$ block matrix. The $(i, j)$th block in $B^{11}$ is expressed as equation 5.7, and $O(n)$ flops are needed to compute these coefficients $a_{ij}$, $b_{kl}^{ij}$ and $z_{kl}^{ij}$.

The computation of the natural gradient is based on theorem 1. Let

$$
\widetilde{\nabla}l_1 = \begin{bmatrix} (\widetilde{\nabla}l_1)_{\boldsymbol{w}_1} \\ \vdots \\ (\widetilde{\nabla}l_1)_{\boldsymbol{w}_m} \\ (\widetilde{\nabla}l_1)_{\boldsymbol{a}} \\ (\widetilde{\nabla}l_1)_{\boldsymbol{b}} \end{bmatrix}
$$

be the partition of $\widetilde{\nabla}l_1$. By theorem 1,

$$
(\widetilde{\nabla}l_1)_{\boldsymbol{w}_1} = \sum_{i=1}^{m} \widetilde{\boldsymbol{B}}_{1i}\frac{\partial l_1}{\partial \boldsymbol{w}_i} - \boldsymbol{Y}_{11}\frac{\partial l_1}{\partial \boldsymbol{a}} - \boldsymbol{Y}_{12}\frac{\partial l_1}{\partial \boldsymbol{b}}.
$$

Applying equations 5.6 and 5.7, we have

$$
(\widetilde{\nabla}l_1)_{\boldsymbol{w}_1} = \sum_{k=1}^{m} \left( a_{1k}\frac{\partial l_1}{\partial \boldsymbol{w}_k} + \widetilde{a}_{1k}\boldsymbol{u}_k \right), \tag{5.8}
$$

where $a_{1k}$ are the coefficients of $\widetilde{\boldsymbol{B}}_{1k}$ in equation 5.7 and

$$
\widetilde{a}_{1k} = -\sum_{i=1}^{m} a_{1i}\boldsymbol{v}_k^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} + \sum_{i,l=1}^{m}\left( b_{kl}^{1i}\boldsymbol{v}_l^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} + z_{kl}^{1i}\boldsymbol{u}_l^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} \right)
$$

$$
-\sum_{l=1}^{m}\left( y_{11}^{kl}\frac{\partial l_1}{\partial \boldsymbol{a}_l} + y_{12}^{kl}\frac{\partial l_1}{\partial \boldsymbol{b}_l} \right).
$$

Therefore, $(\widetilde{\nabla}l_1)_{\boldsymbol{w}_1}$ is a linear combination of $\frac{\partial l_1}{\partial \boldsymbol{w}_i}$ and $\boldsymbol{u}_k$, and it requires $O(n)$ flops to compute $a_{1k}$ and $\widetilde{a}_{1k}$ due to the vector inner product and the vector addition. This is also true for computing other $(\widetilde{\nabla}l_1)_{\boldsymbol{w}_i}$.

Let $\boldsymbol{B}^{22} = [C_{ij}]_{[2\times2]}$. Again applying equation 5.6, we have

$$
(\widetilde{\nabla}l_1)_{\boldsymbol{a}} = \sum_{i=1}^{m} \boldsymbol{Y}_{i1}^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} + C_{11}\frac{\partial l_1}{\partial \boldsymbol{a}} + C_{12}\frac{\partial l_1}{\partial \boldsymbol{b}}, \tag{5.9}
$$

where

$$
\boldsymbol{Y}_{i1}^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} = \left[ \sum_{k=1}^{m} y_{i1}^{k1}\boldsymbol{u}_k^T\frac{\partial l_1}{\partial \boldsymbol{w}_i}, \cdots, \sum_{k=1}^{m} y_{i1}^{km}\boldsymbol{u}_k^T\frac{\partial l_1}{\partial \boldsymbol{w}_i} \right]^T
$$

involving the vector inner product $\boldsymbol{u}_k^T\frac{\partial l_1}{\partial \boldsymbol{w}_i}$. Note that each $C_{ij}$ is of $m \times m$. So, it only needs $O(n)$ flops to compute $(\widetilde{\nabla}l_1)_{\boldsymbol{a}}$. This is also true for $(\widetilde{\nabla}l_1)_{\boldsymbol{b}}$.

**Corollary 1.**   *The time complexity for inverting the Fisher information matrix is $O(n^2)$. The time complexity for computing the natural gradient $\boldsymbol{A}(\boldsymbol{\theta})^{-1}\frac{\partial l_1}{\partial \boldsymbol{\theta}}$ is $O(n)$.*

In summary, it is easy to generalize formulas 5.8 and 5.9 to compute $(\tilde{\nabla} l_1)_{w_i}$ and $(\tilde{\nabla} l_1)_b$. For the multilayer perceptron, the natural gradient can be computed in $O(n)$ flops by these formulas.

## 6 Conclusion

The complexity for computing the Fisher information matrix and the natural gradient is discussed. Based on the efficient scheme to represent the Fisher information matrix, the process of inverting the Fisher information matrix is examined in detail. The structure of the Fisher information matrix and the natural gradient for the multilayer perceptron is elaborated by this analysis. We show that the time complexity for inverting the Fisher information matrix is $O(n^2)$ and the time complexity for computing the natural gradient is $O(n)$.

## Appendix A:

### Proof of Lemma 5

Extend $\{u_1, \cdots, u_m\}$ to $\{u_1, \cdots, u_m, u_{m+1}, \cdots, u_n\}$, a basis in $\mathfrak{R}^{n \times n}$, such that

$$\text{for } j > m, \ u_j \perp \mathcal{L}(u_1, \cdots, u_m) \tag{A.1}$$

$$u_j^T u_k = \delta_{j,k} \text{ (delta notation)}, \quad \text{for } j, k = m+1, \cdots, n, \tag{A.2}$$

where $\mathcal{L}(u_1, \cdots, u_m)$ is the vector space spanned by $\{u_1, \cdots, u_m\}$.

Let $U = [U_1, u_{m+1}, \cdots, u_n]$ and $V = [V_1, u_{m+1}, \cdots, u_n]$. Then noticing that each column of $V_1$ is in $\mathcal{L}(u_1, \cdots, u_m)$, we have $V^T U = I$. Therefore, $V^T = U^{-1}$ and $UV^T = I$, which implies that

$$\sum_{i=1}^m u_i v_i + \sum_{i=m+1}^n u_i u_i^T = I.$$

that is,

$$\Omega_0 = I - \sum_{i=1}^m u_i v_i = \sum_{i=m+1}^n u_i u_i^T. \tag{A.3}$$

Let $A = a_0 \Omega_0 + \sum_{ij} a_{ij} u_i v_j^T$ and $B \in \mathcal{M}_v$. Each column of $B$ is a linear combination of $\{v_1, \cdots, v_m\}$. Let $b = \sum_{i=1}^m c_i v_i$. Then

$$Ab = \left( a_0 \sum_{i=m+1}^n u_i u_i^T + \sum_{i,j=1}^m a_{ij} u_i v_j^T \right) \sum_{i=1}^m c_i v_i$$

$$= \sum_{i,j,k=1}^{m} a_{ij} c_k v_j^T v_k u_i$$

$$= \sum_{i=1}^{m} a^i u_i,$$

where $a^i = \sum_{j,k=1}^{m} a_{ij} c_k v_j^T v_k$ and $O(n)$ flops are needed to compute these coefficients $a^i$. Expressing $AB$ as the following form,

$$AB = \left[ \sum_{i=1}^{m} a_1^i u_i, \cdots, \sum_{i=1}^{m} a_m^i u_i \right],$$

we need $O(n)$ flops to compute the coefficients $a_j^i$.

**Appendix B:** ────────────────────────────────────

**Proof of Lemma 6**

Let $A = [\sum_{i=1}^{m} a_1^i u_i, \cdots, \sum_{i=1}^{m} a_m^i u_i]$ and $B = [\sum_{i=1}^{m} b_1^i v_i, \cdots, \sum_{i=1}^{m} b_m^i v_i]$. Then the (k,l)th element of $B^T A$ is $\sum_{i,j=1}^{m} b_k^i a_l^j v_i^T u_j$, which is computed in $O(n)$ flops.

**Appendix C:** ────────────────────────────────────

**Proof of Lemma 7**

For $A \in \mathcal{M}_u$ and $B \in \Re^{m \times m}$, each column of $AB$ has the form $Ax$ for $x \in \Re^m$. Let $A = [\sum_{i=1}^{m} c_1^i u_i, \cdots, \sum_{i=1}^{m} c_m^i u_i]$ and $x = [x_1, \cdots, x_m]^T$. Then

$$Ax = \sum_{i=1}^{m} \left( \sum_{k=1}^{m} c_k^i x_k \right) u_i = \sum_{i=1}^{m} b^i u_i,$$

and $m^2$ flops are needed to compute these coefficients $b^i$. Hence, when $AB$ is expressed as the following form,

$$AB = \left[ \sum_{i=1}^{m} b_1^i u_i, \cdots, \sum_{i=1}^{m} b_m^i u_i \right],$$

$m^3$ flops are needed to compute the coefficients $b_j^i$.

**Appendix D:** _____

### Proof of Lemma 8

$B_{11}$ has the same structure as that of the Fisher information matrix for the committee machine. $B_{11}^{-1}$ is an $m \times m$ block matrix, and each of its $m^2$ blocks has the following form,

$$\widetilde{a}_0 \Omega_0 + \sum_{l,k=1}^{m} \widetilde{c}_{ij}^{lk} u_l v_k^T,$$

and the complexity for computing these coefficients $\widetilde{a}_0$ and $\widetilde{c}_{ij}^{lk}$ does not depend on $n$.

$B_{12}$ is an $m \times 2$ block matrix with blocks belonging to $\mathcal{M}_v$. By lemma 5, $B_{11}^{-1} B_{12}$ is an $m \times 2$ block matrix, and each of its $2m$ blocks belongs to $\mathcal{M}_u$ with the form of equation 5.5 and the complexity for computing $x_{ij}^{kl}$ is $O(n)$.

**Appendix E:** _____

### Proof of Lemma 9

Applying lemmas 6 and 8, it takes $O(n)$ flops to compute the $2m \times 2m$ matrix $B_{12}^T (B_{11}^{-1} B_{12})$, so $O(n)$ flops are needed to compute $B_{22,1} = B_{22} - B_{12}^T (B_{11}^{-1} B_{12})$.

By lemmas 7 and 8,

$$B_{11}^{-1} B_{12} B_{22,1}^{-1} = \begin{bmatrix} Y_{11} & Y_{12} \\ \vdots & \vdots \\ Y_{m1} & Y_{m2} \end{bmatrix},$$

where $Y_{ij} = [\sum_{k=1}^{m} y_{ij}^{k1} u_k, \cdots, \sum_{k=1}^{m} y_{ij}^{km} u_k]$.

Also by lemma 8,

$$B_{21} B_{11}^{-1} = B_{12}^T B_{11}^{-1} = (B_{11}^{-1} B_{12})^T = \begin{bmatrix} X_{11}^T & \cdots & X_{m1}^T \\ X_{12}^T & \cdots & X_{m2}^T \end{bmatrix}.$$

Hence,

$$B_{11}^{-1} B_{12} B_{22,1}^{-1} B_{21} B_{11}^{-1} = [Z_{ij}]_{[m \times m]},$$

where $Z_{ij} = Y_{i1} X_{j1}^T + Y_{i2} X_{j2}^T = \sum_{k,l=1}^{m} z_{ij}^{kl} u_k u_l^T$ and $O(n)$ flops are needed to compute the coefficients $z_{ij}^{kl}$.

## References

Amari, S. (1997). Neural learning in structured parameter spaces—Natural Riemannian gradient. In M. C. Mozer, M. I. Jordan, & T. Petsche (Eds.), *Advances in neural information processing systems, 9* (pp. 127–133) Cambridge, MA: MIT Press.

Amari, S. (1998). Natural gradient works efficiently in learning. *Neural Computation, 10,* 251–276.

Orr, G. B., & Leen, T. K. (1997). Using curvature information for fast stochastic search. In M. C. Mozer, M. I. Jordan & T. Petsche (Eds.), *Advances in neural information processing systems, 9.* Cambridge, MA: MIT Press.

Saad, D., & Solla, S. A. (1995). On-line learning in soft committee machines. *Physical Review E, 52,* 4225–4243.

Stewart, G. W. (1973). *Introduction to matrix computations.* New York: Academic Press.

Yang, H. H., & Amari, S. (1997a). Natural gradient descent for training multilayer perceptrons. Unpublished manuscript.

Yang, H. H., & Amari, S. (1997b). Training multi-layer perceptrons by natural gradient descent. In *ICONIP'97 Proceedings*, New Zealand.

Yang, H. H., & Amari, S. (1998). The efficiency and the robustness of natural gradient descent learning rule. In M. I. Jordan, M. J. Kearns, & S. A. Solla (Eds.), *Advances in neural information processing systems, 10.* Cambridge, MA: MIT Press.

Yang, H. H., Murata, N., & Amari, S. (1998). Statistical inference: Learning in artificial neural networks. *Trends in Cognitive Sciences, 2*(1), 4–10.