

## 多層パーセプトロンにおける自然勾配学習

理化学研究所 脳科学総合研究センター 甘利 俊一

### [要 旨]

多層パーセプトロンは、多数のパラメータを含む非線形の関数を用いてその入出力関係を表現する。学習とは、実現したい入出力関係が外部より与えられたときに、その例題のみに基づいてパーセプトロンの可変パラメータを調整し、これによって例題の含む“一般的”な仕組みを得てこの関数を近似的に実現することである。特に、オンライン学習は、例題が一つずつ逐次的に与えられたときに、これを記憶することなくパラメータ学習を行なう。

学習アルゴリズムの典型的なものは、誤差関数の確率降下法に基づくもので、誤差逆伝播法の名前で知られている。ところが、これは学習の収束がきわめて遅いところに問題がある。

近年の統計物理学的手法によるオンライン学習の研究は、この収束の遅さがプラトーと呼ばれる対称性に起因する領域に状態が捉わるからであること、さらにプラトーは普遍的に存在することを明らかにした。

ではプラトーから脱出する方策として何がよいであろうか。本研究は、パーセプトロンのパラメータ空間はユークリッド空間ではなくてリーマン空間の構造を持つことを明らかにする。リーマン空間における勾配（グラディエント）は反変ベクトルで表現されるべきであり、これが真の最急方向を与える。この自然のリーマン構造を用いた学習法は、局所的に漸近最適であるのみならず、プラトーに捉われないものであることを示す。

# Natural Gradient Eliminates Plateaus in Multilayer Perceptrons Learning

Shun-ichi Amari  
RIKEN Brain Science Institute  
amari@brain.riken.go.jp

## Abstract

The present memo studies the topological and metrical structures of the parameter space of neural networks, in particular multilayer perceptrons. It is important to study such geometrical structures for evaluating learning abilities of various algorithm. We further propose the natural Riemannian gradient learning algorithm which is powerful not only in the asymptotic regime but also in the intermediate regime.

## 1 Multilayer perceptrons

Let us consider the following model of multilayer perceptrons,

$$y = \sum_{\alpha=1}^m v_{\alpha} \varphi(\mathbf{w}_{\alpha} \cdot \mathbf{x}) + \xi, \quad (1)$$

where  $\mathbf{x}$  is an  $n$ -dimensional input subject to the Gaussian distribution  $N(O, I)$ ,  $I$  being the  $n \times n$  identity matrix;  $y$  is an output from a linear output unit;  $\mathbf{w}_{\alpha}$  is an  $n$ -dimensional connection weight vector from the input to the  $\alpha$ -th hidden unit,  $\alpha = 1, \dots, m$ ,  $m$  being the number of hidden units;  $v_{\alpha}$  is the connection weight from the  $\alpha$ -th hidden unit to the output unit; and  $\xi$  is a random noise subject to  $N(0, \sigma^2)$ . The function  $\varphi$  is a sigmoidal function, and we use

$$\varphi(u) = \sqrt{\frac{2}{\pi}} \int_0^u \exp\left\{-\frac{z^2}{2}\right\} dz \quad (2)$$

in the present memo in order to obtain an explicit analytical form of the Fisher information matrix. When bias terms exist, the output of the  $\alpha$ -th unit is written as

$$\varphi(\mathbf{w}_{\alpha} \cdot \mathbf{x} + \mathbf{b}_{\alpha}), \quad (3)$$

but we neglect  $\mathbf{b}$  in this memo for simplicity's sake. The terms  $\mathbf{b}_{\alpha}$  can be analyzed in a similar way without difficulty. When all  $v_{\alpha}$  are fixed to be equal to 1, we have a committee machine,

$$y = \sum \varphi(\mathbf{w}_{\alpha} \cdot \mathbf{x}) + \xi, \quad (4)$$

as a special case.

The parameters  $\{\mathbf{w}_1, \dots, \mathbf{w}_{\alpha}; \mathbf{v}\}$  can be summarized into a single  $m(n+1)$ -dimensional vector  $\theta$ , which plays the role of a (local) coordinate system in the space  $S$  of the multilayer perceptrons.

A perceptron having parameter  $\theta$  gives the conditional probability of output  $y$  conditioned on input  $\mathbf{x}$ ,

$$p(y|\mathbf{x}, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2\sigma^2} \{y - f(\mathbf{x}, \theta)\}^2 \right], \quad (5)$$

$$f(\mathbf{x}, \theta) = \sum v_\alpha \varphi(\mathbf{w}_\alpha \cdot \mathbf{x}) \quad (6)$$

is the mean value of  $y$  given input  $\mathbf{x}$ . Its logarithm is

$$l(z|\mathbf{x}, \theta) = -\frac{1}{2\sigma^2} \{y - f(\mathbf{x}, \theta)\}^2 - \log(\sqrt{2\pi}\sigma). \quad (7)$$

This can be regarded as the negative of the square of an error when  $y$  is a target value, except for a scale and a constant term. The space  $S$  of perceptrons is identified with the set of all the conditional probability distributions of the form (5).

## 2 Global Structure of $S$

When

$$p(y|\mathbf{x}; \theta_1) = p(y|\mathbf{x}; \theta_2) \quad (8)$$

holds, two perceptrons with parameters  $\theta_1$  and  $\theta_2$  have the same input-output relation. We define that two perceptrons (or two points)  $\theta_1$  and  $\theta_2$  are equivalent,  $\theta_1 \approx \theta_2$  in this case. When  $\theta$  has no equivalent points other than itself, it is called a proper point. Otherwise, it is a singular point. From the behavioral point of view, we need to consider the quotient space  $\tilde{S} = S/R$ , where  $R$  is the equivalence relation just introduced.

In order to study the topological structure of  $\tilde{S}$ , we study when two points are equivalent. To this end, we introduce the following group acting on  $S$  (Chen, Liu and Hecht-Nielsen). Let  $\Pi$  be a permutation of indices  $\{1, 2, \dots, m\}$  which maps  $i$  to  $\Pi(i)$ . Let  $\Pi$  act on  $\theta = \{\mathbf{w}_1, \dots, \mathbf{w}_m; v_1, \dots, v_m\}$  as

$$\Pi \{\mathbf{w}_1, \dots, \mathbf{w}_m; v\} = \{\mathbf{w}_{\Pi(1)}, \dots, \mathbf{w}_{\Pi(m)}; v_{\Pi(1)}, \dots, v_{\Pi(m)}\}. \quad (9)$$

Let  $C_\alpha (\alpha = 1, \dots, m)$  be the change of signs of  $\mathbf{w}_\alpha$  and  $v_\alpha$ :

$$C_\alpha \{\mathbf{w}_1, \dots, \mathbf{w}_m; v_1, \dots, v_m\} = \{\mathbf{w}_1, \dots, -\mathbf{w}_\alpha, \dots, \mathbf{w}_m; v_1, \dots, -v_\alpha, \dots, v_m\}. \quad (10)$$

**Theorem 1 (Chen, Liu and Hecht-Nielsen).** Two points are equivalent when they are transformed to each other by an element of the group composed of  $Gr = \{\Pi, C_\alpha\}$ .

The proof is trivial since the form (1) is invariant under  $\Pi$  and  $C_\alpha$  because of  $\varphi(-u) = -\varphi(u)$ .

Therefore, we need to study the residual set  $S/Gr$ . However, when  $\theta$  is invariant under some elements of  $Gr$ , further singular points are found. We have the following three types:

1. When  $v_\alpha = 0$ , two points are equivalent when they differ only by  $\mathbf{w}_\alpha$ .
2. When  $\mathbf{w}_\alpha = 0$ , two points are equivalent when they differ only by  $v_\alpha$ .
3. When  $\mathbf{w}_\alpha = \pm \mathbf{w}_\beta$  holds, two points are equivalent when the values of  $v_\alpha \pm v_\beta$  are equal.

**Theorem 2 (Sussman).** Two points are equivalent when, and only when, they are transformed by an element of  $Gr$  or they are connected by one of the above relations.

We call the subspaces defined by

$$\begin{aligned} 1) \quad & v_\alpha = 0, \\ 2) \quad & \mathbf{w}_\alpha = 0, \\ 3) \quad & \mathbf{w}_\alpha = \pm \mathbf{w}_\beta \end{aligned} \tag{11}$$

critical subspaces. A critical subspace is decomposed into a family of lower dimensional subspaces in which all the points are equivalent. For example,  $v_\alpha = 0$  for a fixed  $\alpha$  defines a hyperplane  $V_\alpha$  in  $S$ . The set of the points which differ only by  $\mathbf{w}_\alpha$  is an  $n$ -dimensional subspace in it, and  $V_\alpha$  is composed of such  $n$ -dimensional subspaces in which all the points are equivalent. Similarly,  $\mathbf{w}_\alpha = \mathbf{w}_\beta$  defines an  $nm$ -dimensional subspace in  $S$ . It is decomposed into a family of lines defined by  $v_\alpha + v_\beta = c$ , on which all the points are equivalent.

It should be noted that critical subspaces may intersect. The space of multilayer perceptrons is the residue class by these equivalence relations. It is an  $m(n+1)$ -dimensional manifold except for the critical submanifolds in which  $\tilde{S}$  has lower-dimensional structures. The following simple example shows such a structure: Let  $S$  be a Euclidean 3-space with coordinates  $(x, y, z)$ . Let  $x = 0$  be a critical subspace, and define an equivalence relation that  $(0, y, z)$  is equivalent to  $(0, y', z)$ . An equivalence class is a line defined by  $x = 0$  and  $z = c$ , and such lines form the critical surface  $x = 0$ . In this example,  $S$  is divided into two parts by  $x = 0$ , and the degeneration of dimension reduction takes place on this surface.

### 3 Structural Degeneration

In the critical subspaces, a multilayer perceptron degenerates into a simpler structure. When  $v_\alpha = 0$  or  $\mathbf{w}_\alpha = 0$  holds, the  $\alpha$ -th hidden unit is ineffective and we can remove it, having a perceptron with  $m - 1$  hidden units. Similarly, when  $\mathbf{w}_\alpha = \mathbf{w}_\beta$  ( $\mathbf{w}_\alpha = -\mathbf{w}_\beta$ ) holds, the outputs of the  $\alpha$ -th and  $\beta$ -th elements are always the same (the same with opposite signs). Therefore, the two units can be merged into one unit where  $v_\alpha + v_\beta$  ( $v_\alpha - v_\beta$ ) is the connection weight of the new unit to the output unit.

Further structural degenerations take place at the intersections of critical subspaces.

### 4 Learning of Multilayer Perceptrons

A multilayer perceptron is trained by examples  $\{\mathbf{x}_1, \mathbf{y}_1^*\}, \dots, \{\mathbf{x}_t, \mathbf{y}_t^*\}$ , where  $\mathbf{y}_t^*$  is the output of the teacher when  $\mathbf{x}_t$  is given. The teacher signal  $\mathbf{y}_t^*$  is generated by

$$\mathbf{y}_t^* = f^*(\mathbf{x}_t) + \xi_t. \tag{12}$$

The conventional on-line (or off-line) learning method modifies the current parameter  $\theta_t$  by using the gradient of the loss function such that

$$\theta_{t+1} = \theta_t - \eta \frac{\partial l^*(\mathbf{x}_t, \mathbf{y}_t^*; \theta_t)}{\partial \theta}. \tag{13}$$

Here  $\mathbf{y}_t^*$  is the teacher output for  $\mathbf{x}_t$  and

$$l^*(\mathbf{x}_t, \mathbf{y}_t^*; \theta_t) = \frac{1}{2} \{y_t^* - f(\mathbf{x}_t; \theta_t)\}^2. \tag{14}$$

This is the conventional backpropagation method. Such a stochastic descent method for analog multilayer perceptron was proposed by Amari (1967), where he suggested to use a positive definite matrix  $C$  to obtain

$$\theta_{t+1} = \theta_t - \eta C \frac{\partial l(\mathbf{x}_t, \mathbf{y}_t^*; \theta_t)}{\partial \theta}. \quad (15)$$

After the monumental work of Rumelhart et al., this method is called the backpropagation learning algorithm, and a lot of acceleration methods have been studied. They are, for example, the momentum method, conjugate gradient, Newton's method, etc. Adaptive learning rate  $\eta$  has also been studied (Amari, 1967 ; Barkai and Sompolinski, Murata et al).

It is known that the backpropagation method is extremely slow because of the existence of plateaus. A plateau is a saddle point of the learning dynamics, and trajectories of dynamics are first attracted by such plateaus. It takes long learning time to get rid of a plateau. Therefore, a typical learning curve takes the form shown in Fig. 2. Where do plateaus exist in  $S$ ? We show that plateaus exist in critical subspaces.

Let us consider a critical subspace  $W_{\alpha\beta}$  defined by  $\mathbf{w}_\alpha = \mathbf{w}_\beta$  ( $\alpha \neq \beta$ ). This corresponds to networks of  $m - 1$  hidden units where the  $\alpha$ - and  $\beta$ -th units are merged. Given a target function  $f^*(\mathbf{x})$  provided by the teacher, let us consider its optimal approximator in  $W_{\alpha\beta}$ , that is, the best approximator of  $f^*(\mathbf{x})$  by a network with  $m - 1$  hidden units. Let it be  $\theta^*$ . From

$$\frac{\partial E[l(\mathbf{x}, \mathbf{y}^*; \theta^*)]}{\partial \theta^*} = 0, \quad (16)$$

where

$$l^* = \frac{1}{2} \{y^* - f(\mathbf{x}, \theta)\}^2. \quad (17)$$

We then have

$$\frac{\partial E[l]}{\partial \mathbf{w}_\gamma} = -E \left[ \{y^* - f(\mathbf{x}, \theta^*)\} \frac{\partial f(\mathbf{x}, \theta^*)}{\partial \mathbf{w}_\gamma} \right] = 0. \quad (18)$$

This gives

$$E \left[ \{y^* - f(\mathbf{x}, \theta^*)\} \varphi'(\mathbf{w}_\gamma^* \cdot \mathbf{x}) \right] = 0 \quad (19)$$

for all  $\gamma$ .

Now we show the following theorem.

**Theorem 3.** The best approximation  $\theta^* \in W_{\alpha\beta}$  is a critical point forming a plateau.

**Proof.** Let us calculate  $\Delta l^*$  corresponding to a change of  $\mathbf{w}_\alpha = \mathbf{w}_\beta$  to  $\mathbf{w}_\alpha + \boldsymbol{\varepsilon}$  and  $\mathbf{w}_\beta - \boldsymbol{\varepsilon}$  for a small vector  $\boldsymbol{\varepsilon}$  by which  $\theta$  escapes from  $W_{\alpha\beta}$ . It is easily shown that

$$\begin{aligned} \Delta l^* &= E \left[ \{y^* - f(\mathbf{x}, \theta^*)\} \varphi'(\mathbf{w}_\alpha \cdot \mathbf{x}) (v_\alpha - v_\beta) \boldsymbol{\varepsilon} \cdot \mathbf{x} \right] \\ &= 0. \end{aligned} \quad (20)$$

Therefore, the gradient of  $l^*$  in directions from  $\theta^*$  toward outside  $W_{\alpha\beta}$  is 0 as well as those inside  $W_{\alpha\beta}$ . Hence,  $\theta^*$  is a critical point. We next show that there exist trajectories staying in  $W_{\alpha\beta}$  and converging to  $\theta^*$ . To show this, we calculate  $\Delta l^*$  at any point  $\theta \in W_{\alpha\beta}$  corresponding to a change of  $\mathbf{w}_\alpha$  by  $v_\beta \boldsymbol{\varepsilon}$  and a change of  $\mathbf{w}_\beta$  ( $= \mathbf{w}_\alpha$ ) by  $v_\alpha \boldsymbol{\varepsilon}$ . The corresponding changes  $\Delta \mathbf{w}_\alpha$  and  $\Delta \mathbf{w}_\beta$  by learning are

$$\Delta \mathbf{w}_\alpha = -\eta \frac{\partial l^*}{\partial \mathbf{w}_\alpha} = -\eta e \varphi'(\mathbf{w}_\alpha \cdot \mathbf{x}) v_\alpha \mathbf{x}, \quad (21)$$

$$\Delta \mathbf{w}_\beta = -\eta \frac{\partial l^*}{\partial \mathbf{w}_\beta} = -\eta e \varphi'(\mathbf{w}_\alpha \cdot \mathbf{x}) v_\beta \mathbf{x}. \quad (22)$$

This shows that

$$\Delta \mathbf{w}_\alpha \cdot (v_\beta \boldsymbol{\varepsilon}) - \Delta \mathbf{w}_\beta \cdot (v_\alpha \boldsymbol{\varepsilon}) = 0 \quad (23)$$

Therefore, at positions where  $v_\alpha = v_\beta$  hold, the gradient flow is always in  $W_{\alpha\beta}$ . This proves that  $\boldsymbol{\theta}^*$  is a plateau.

## 5 Fisher Information Metric and Riemannian Structure

The Kullback-Leibler divergence

$$\begin{aligned} D[\boldsymbol{\theta} : \boldsymbol{\theta}'] &= E \left[ \log \frac{p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}')} \right] \\ &= \int p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \log \frac{p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}')} p(\mathbf{x}) d\mathbf{x} d\mathbf{y} \end{aligned} \quad (24)$$

is a divergence measure between two points  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}'$  in  $S$ . This measure is based on the stochastic behaviors of two networks, so that it is equal to 0 when the two networks  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}'$  are equivalent.

When  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}' = \boldsymbol{\theta} + d\boldsymbol{\theta}$  are infinitesimally close to each other, we have

$$\begin{aligned} D(\boldsymbol{\theta} : \boldsymbol{\theta} + d\boldsymbol{\theta}) &= \frac{1}{2} d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \frac{1}{2} \sum g_{ij}(\boldsymbol{\theta}) d\theta_i d\theta_j, \end{aligned} \quad (25)$$

where  $G(\boldsymbol{\theta}) = (g_{ij}(\boldsymbol{\theta}))$  is a matrix defined by

$$\begin{aligned} g_{ij}(\boldsymbol{\theta}) &= E \left[ \frac{\partial l(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial l(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j} \right] \\ &= -E \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \right]. \end{aligned} \quad (26)$$

This is the Fisher information matrix.

The Fisher information defines a Riemannian metric in the tangent space  $T_{\boldsymbol{\theta}}$  at each  $\boldsymbol{\theta}$ . Given a vector  $d\boldsymbol{\theta}$  which represents a small deviation of  $\boldsymbol{\theta}$  at  $\boldsymbol{\theta}$ , its squared length is defined by

$$\langle d\boldsymbol{\theta}, d\boldsymbol{\theta} \rangle_{\boldsymbol{\theta}} = d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (27)$$

Two vectors  $d\boldsymbol{\theta}$  and  $d\boldsymbol{\theta}'$  are said to be orthogonal when

$$\langle d\boldsymbol{\theta}, d\boldsymbol{\theta}' \rangle_{\boldsymbol{\theta}} = d\boldsymbol{\theta}^T G(\boldsymbol{\theta}) d\boldsymbol{\theta}' = 0. \quad (28)$$

When  $G(\boldsymbol{\theta}) = I$ , the identity matrix, at any  $\boldsymbol{\theta}$ , the space reduces to a Euclidean space and  $\langle , \rangle$  is the ordinary inner product in the Euclidean space. A space is said to be Riemannian, when  $G(\boldsymbol{\theta})$  is defined at each  $\boldsymbol{\theta}$ . The space of multilayer perceptrons is intrinsically Riemannian, not reducing to a Euclidean one.

The Riemannian metric represents the topological structure of  $S$  very well. In a critical subspace, for example  $W_{\alpha\beta}$ , two points  $(\mathbf{w}_1, \dots, \mathbf{w}_m; \mathbf{v})$  and  $(\mathbf{w}_1, \dots, \mathbf{w}_m; \mathbf{v}')$  are equivalent

when  $\mathbf{w}_\alpha = \mathbf{w}_\beta$  and  $v_\alpha + v_\beta = v'_\alpha + v'_\beta$ . A small deviation  $d\theta$  in the equivalent direction  $d\theta$  in this equivalence class is of length 0,

$$d\theta^T G(\theta) d\theta = 0, \quad (29)$$

implying that  $G(\theta)$  is degenerate in a critical subspace. Therefore, when  $\theta$  is close to a critical subspace, it is ill conditioned. Inside a critical subspace, we can define  $G(\theta)$  in this subspace which corresponds to the space of a reduced architecture of networks. It is regular in the subspace except for critical subspaces inside it.

**Theorem 4 (Fukumizu).** The Fisher information is regular except for critical subspaces.

## 6 Natural Gradient Learning Algorithm

The natural gradient learning algorithm (Amari, 1998) updates the current  $\theta_t$  by

$$\theta_{t+1} = \theta_t - \eta_t G^{-1}(\theta_t) \frac{\partial l(y_t, \mathbf{x}_t; \theta_t)}{\partial \theta} \quad (30)$$

where  $\eta_t$  is a learning constant which may depend on  $t$ . The term

$$\tilde{\nabla} l = G^{-1}(\theta_t) \frac{\partial l}{\partial \theta} \quad (31)$$

gives the true steepest direction of function  $l$  defined in a Riemannian manifold. The inverse  $G^{-1}$  does not exist in critical subspaces. Therefore, one should use the reduced  $G$  in the submanifold.

The natural gradient learning method has the following remarkable property.

**Theorem 5 (Amari, Opper).** When the learning rate is  $\eta_t = 1/t$ , the estimator  $\theta_t$  obtained by the natural gradient learning algorithm is Fisher efficient,

$$E \left[ (\theta_t - \theta)^T (\theta_t - \theta) \right] = \frac{1}{t} G^{-1}(\theta), \quad (32)$$

that is, its asymptotic covariance is the same as the best batch estimator.

A much more important property was suggested by Amari (1998) that the natural gradient learning algorithm eliminates plateaus. We can elucidate how plateaus are eliminated by natural gradient descent learning by using the committee machine.

## References

- [1] Amari, S. (1998), "Natural gradient works efficiently in learning", *Neural Computation*, 10, 251-276.
- [2] Amari, S. (1998), "Natural gradient for over- and under-complete bases in ICA, submitted.