
Neural Learning in Structured Parameter Spaces — Natural Riemannian Gradient

Shun-ichi Amari

RIKEN Frontier Research Program, RIKEN,
Hirosawa 2-1, Wako-shi 351-01, Japan
amari@zoo.riken.go.jp

Abstract

The parameter space of neural networks has a Riemannian metric structure. The natural Riemannian gradient should be used instead of the conventional gradient, since the former denotes the true steepest descent direction of a loss function in the Riemannian space. The behavior of the stochastic gradient learning algorithm is much more effective if the natural gradient is used. The present paper studies the information-geometrical structure of perceptrons and other networks, and prove that the on-line learning method based on the natural gradient is asymptotically as efficient as the optimal batch algorithm. Adaptive modification of the learning constant is proposed and analyzed in terms of the Riemannian measure and is shown to be efficient. The natural gradient is finally applied to blind separation of mixed independent signal sources.

1 Introduction

Neural learning takes place in the parameter space of modifiable synaptic weights of a neural network. The role of each parameter is different in the neural network so that the parameter space is structured in this sense. The Riemannian structure which represents a local distance measure is introduced in the parameter space by information geometry (Amari, 1985).

On-line learning is mostly based on the stochastic gradient descent method, where the current weight vector is modified in the gradient direction of a loss function. However, the ordinary gradient does not represent the steepest direction of a loss function in the Riemannian space. A geometrical modification is necessary, and it is called the natural Riemannian gradient. The present paper studies the remarkable effects of using the natural Riemannian gradient in neural learning.

We first studies the asymptotic behavior of on-line learning (Oppor, NIPS'95 Workshop). Batch learning uses all the examples at any time to obtain the optimal weight vector, whereas on-line learning uses an example once when it is observed. Hence, in general, the target weight vector is estimated more accurately in the case of batch learning. However, we prove that, when the Riemannian gradient is used, on-line learning is asymptotically as efficient as optimal batch learning.

On-line learning is useful when the target vector fluctuates slowly (Amari, 1967). In this case, we need to modify a learning constant η_t depending on how far the current weight vector is located from the target function. We show an algorithm adaptive changes in the learning constant based on the Riemannian criterion and prove that it gives asymptotically optimal behavior. This is a generalization of the idea of Sompolinsky et al. [1995].

We then answer the question what is the Riemannian structure to be introduced in the parameter space of synaptic weights. We answer this problem from the point of view of information geometry (Amari [1985, 1995], Amari et al [1992]). The explicit form of the Riemannian metric and its inverse matrix are given in the case of simple perceptrons.

We finally show how the Riemannian gradient is applied to blind separation of mixed independent signal sources. Here, the mixing matrix is unknown so that the parameter space is the space of matrices. The Riemannian structure is introduced in it. The natural Riemannian gradient is computationaly much simpler and more effective than the conventional gradient.

2 Stochastic Gradient Descent and On-Line Learning

Let us consider a neural network which is specified by a vector parameter $\mathbf{w} = (w_1, \dots, w_n) \in \mathbf{R}^n$. The parameter \mathbf{w} is composed of modifiable connection weights and thresholds. Let us denote by $l(\mathbf{x}, \mathbf{w})$ a loss when input signal \mathbf{x} is processed by a network having parameter \mathbf{w} . In the case of multilayer perceptrons, a desired output or teacher signal \mathbf{y} is associated with \mathbf{x} , and a typical loss is given

$$l(\mathbf{x}, \mathbf{y}, \mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{f}(\mathbf{x}, \mathbf{w})\|^2, \quad (1)$$

where $\mathbf{z} = \mathbf{f}(\mathbf{x}, \mathbf{w})$ is the output from the network.

When input \mathbf{x} , or input-output training pair (\mathbf{x}, \mathbf{y}) , is generated from a fixed probability distribution, the expected loss $L(\mathbf{w})$ of the network specified by \mathbf{w} is

$$L(\mathbf{w}) = E[l(\mathbf{x}, \mathbf{y}; \mathbf{w})], \quad (2)$$

where E denotes the expectation. A neural network is trained by using training examples $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots$ to obtain the optimal network parameter \mathbf{w}^* that minimizes $L(\mathbf{w})$. If $L(\mathbf{w})$ is known, the gradient method is described by

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla L(\mathbf{w}_t), \quad t = 1, 2, \dots$$

where η_t is a learning constant depending on t and $\nabla L = \partial L / \partial \mathbf{w}$. Usually $L(\mathbf{w})$ is unknown. The stochastic gradient learning method

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla l(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathbf{w}_t) \quad (3)$$

was proposed by an old paper (Amari [1967]). This method has become popular since Rumelhart et al. [1986] rediscovered it. It is expected that, when η_t converges to 0 in a certain manner, the above \mathbf{w}_t converges to \mathbf{w}^* . The dynamical behavior of

(3) was studied by Amari [1967], Heskes and Kappen [1992] and many others when η_t is a constant.

It was also shown in Amari [1967] that

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t G^{-1} \nabla l(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{w}_t) \quad (4)$$

works well for any positive-definite matrix, in particular for the metric G . Geometrically speaking $\partial l / \partial \mathbf{w}$ is a covariant vector while $\Delta \mathbf{w}_t = \mathbf{w}_{t+1} - \mathbf{w}_t$ is a contravariant vector. Therefore, it is natural to use a (contravariant) metric tensor G^{-1} to convert the covariant gradient into the contravariant form

$$\tilde{\nabla} l = G^{-1} \frac{\partial l}{\partial \mathbf{w}} = \left(\sum_j g^{ij} \frac{\partial}{\partial w_j}(\mathbf{w}) \right), \quad (5)$$

where $G^{-1} = (g^{ij})$ is the inverse matrix of $G = (g_{ij})$. The present paper studies how the matrix tensor matrix G should be defined in neural learning and how effective is the new gradient learning rule

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \tilde{\nabla} l(\mathbf{x}_t, \mathbf{y}_t, \mathbf{w}_t). \quad (6)$$

3 Gradient in Riemannian spaces

Let $S = \{\mathbf{w}\}$ be the parameter space and let $l(\mathbf{w})$ be a function defined on S . When S is a Euclidean space and \mathbf{w} is an orthonormal coordinate system, the squared length of a small incremental vector $d\mathbf{w}$ connecting \mathbf{w} and $\mathbf{w} + d\mathbf{w}$ is given by

$$|d\mathbf{w}|^2 = \sum_{i=1}^n (dw_i)^2. \quad (7)$$

However, when the coordinate system is non-orthonormal or the space S is Riemannian, the squared length is given by a quadratic form

$$|d\mathbf{w}|^2 = \sum_{i,j} g_{ij}(\mathbf{w}) dw_i dw_j = \mathbf{w}' G \mathbf{w}. \quad (8)$$

Here, the matrix $G = (g_{ij})$ depends in general on \mathbf{w} and is called the metric tensor. It reduces to the identity matrix I in the Euclidean orthonormal case. It will be shown soon that the parameter space S of neural networks has Riemannian structure (see Amari et al. [1992], Amari [1995], etc.).

The steepest descent direction of a function $l(\mathbf{w})$ at \mathbf{w} is defined by a vector $d\mathbf{w}$ that minimize $l(\mathbf{w} + d\mathbf{w})$ under the constraint $|d\mathbf{w}|^2 = \varepsilon^2$ (see eq.8) for a sufficiently small constant ε .

Lemma 1. The steepest descent direction of $l(\mathbf{w})$ in a Riemannian space is given by

$$-\tilde{\nabla} l(\mathbf{w}) = -G^{-1}(\mathbf{w}) \nabla l(\mathbf{w}).$$

We call

$$\tilde{\nabla} l(\mathbf{w}) = G^{-1}(\mathbf{w}) \nabla l(\mathbf{w})$$

the natural gradient of $l(\mathbf{w})$ in the Riemannian space. It shows the steepest descent direction of l , and is nothing but the contravariant form of ∇l in the tensor notation. When the space is Euclidean and the coordinate system is orthonormal, G is the unit matrix I so that $\tilde{\nabla} l = \nabla l$.

4 Natural gradient gives efficient on-line learning

Let us begin with the simplest case of noisy multilayer analog perceptrons. Given input \mathbf{x} , the network emits output $\mathbf{z} = \mathbf{f}(\mathbf{x}, \mathbf{w}) + \mathbf{n}$, where \mathbf{f} is a differentiable deterministic function of the multilayer perceptron with parameter \mathbf{w} and \mathbf{n} is a noise subject to the normal distribution $N(0, I)$. The probability density of an input-output pair (\mathbf{x}, \mathbf{z}) is given by

$$p(\mathbf{x}, \mathbf{z}; \mathbf{w}) = q(\mathbf{x})p(\mathbf{z}|\mathbf{x}; \mathbf{w}),$$

where $q(\mathbf{x})$ is the probability distribution of input \mathbf{x} , and

$$p(\mathbf{z}|\mathbf{x}; \mathbf{w}) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \|\mathbf{z} - \mathbf{f}(\mathbf{x}, \mathbf{w})\|^2 \right\}.$$

The squared error loss function (1) can be written as

$$l(\mathbf{x}, \mathbf{z}, \mathbf{w}) = -\log p(\mathbf{x}, \mathbf{z}; \mathbf{w}) + \log q(\mathbf{x}) - \log \sqrt{2\pi}.$$

Hence, minimizing the loss is equivalent to maximizing the likelihood function $p(\mathbf{x}, \mathbf{z}; \mathbf{w})$.

Let $D_T = \{(\mathbf{x}_1, \mathbf{z}_1), \dots, (\mathbf{x}_T, \mathbf{z}_T)\}$ be T independent input-output examples generated by the network having the parameter \mathbf{w}^* . Then, maximum likelihood estimator $\hat{\mathbf{w}}_T$ minimizes the log loss $l(\mathbf{x}, \mathbf{z}; \mathbf{w}) = -\log p(\mathbf{x}, \mathbf{z}; \mathbf{w})$ over the training data D_T , that is, it minimizes the training error

$$E_{\text{train}}(\mathbf{w}) = \frac{1}{T} \sum_{t=1}^T l(\mathbf{x}_t, \mathbf{z}_t; \mathbf{w}). \quad (9)$$

The maximum likelihood estimator is efficient (or Fisher-efficient), implying that it is the best consistent estimator satisfying the Cramér-Rao bound asymptotically,

$$\lim_{T \rightarrow \infty} TE[(\hat{\mathbf{w}}_T - \mathbf{w}^*)(\hat{\mathbf{w}}_T - \mathbf{w}^*)'] = G^{-1}, \quad (10)$$

where G^{-1} is the inverse of the Fisher information matrix $G = (g_{ij})$ defined by

$$g_{ij} = E \left[\frac{\partial \log p(\mathbf{x}, \mathbf{z}; \mathbf{w})}{\partial w_i} \frac{\partial \log p(\mathbf{x}, \mathbf{z}; \mathbf{w})}{\partial w_j} \right] \quad (11)$$

in the component form. Information geometry (Amari, 1985) proves that the Fisher information G is the only invariant metric to be introduced in the space $S = \{\mathbf{w}\}$ of the parameters of probability distributions.

Examples $(\mathbf{x}_1, \mathbf{z}_1), (\mathbf{x}_2, \mathbf{z}_2) \dots$ are given one at a time in the case of on-line learning. Let $\tilde{\mathbf{w}}_t$ be the estimated value at time t . At the next time $t+1$, the estimator $\tilde{\mathbf{w}}_t$ is modified to give a new estimator $\tilde{\mathbf{w}}_{t+1}$ based on the observation $(\mathbf{x}_{t+1}, \mathbf{z}_{t+1})$. The old observations $(\mathbf{x}_1, \mathbf{z}_1), \dots, (\mathbf{x}_t, \mathbf{z}_t)$ cannot be reused to obtain $\tilde{\mathbf{w}}_{t+1}$, so that the learning rule is written as $\tilde{\mathbf{w}}_{t+1} = \mathbf{n}(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}, \tilde{\mathbf{w}}_t)$. The process $\{\tilde{\mathbf{w}}_t\}$ is hence Markovian. Whatever a learning rule \mathbf{n} we choose, the behavior of the estimator $\tilde{\mathbf{w}}_t$ is never better than that of the optimal batch estimator $\hat{\mathbf{w}}_t$ because of this restriction. The conventional on-line learning rule is given by the following gradient form $\tilde{\mathbf{w}}_{t+1} = \tilde{\mathbf{w}}_t - \eta_t \nabla l(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}; \tilde{\mathbf{w}}_t)$. When η_t satisfies a certain condition, say $\eta_t = c/t$, the stochastic approximation guarantees that $\tilde{\mathbf{w}}_t$ is a consistent estimator converging to \mathbf{w}^* . However, it is not efficient in general.

There arises a question if there exists an on-line learning rule that gives an efficient estimator. If it exists, the asymptotic behavior of on-line learning is equivalent to

that of the batch estimation method. The present paper answers the question by giving an efficient on-line learning rule

$$\tilde{\mathbf{w}}_{t+1} = \tilde{\mathbf{w}}_t - \frac{1}{t} \tilde{\nabla} l(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}; \tilde{\mathbf{w}}_t). \quad (12)$$

Theorem 1. The natural gradient on-line learning rule gives an Fisher-efficient estimator, that is,

$$\tilde{V}_t = E[(\tilde{\mathbf{w}}_t - \mathbf{w}^*)(\tilde{\mathbf{w}}_t - \mathbf{w}^*)'] \approx \frac{1}{t} G^{-1}(\mathbf{w}^*). \quad (13)$$

5 Adaptive modification of learning constant

We have proved that $\eta_t = 1/t$ with the coefficient matrix G^{-1} is the asymptotically best choice for on-line learning. However, when the target parameter \mathbf{w}^* is not fixed but fluctuating or changes suddenly, this choice is not good, because the learning system cannot follow the change if η_t is too small. It was proposed in Amari [1967] to choose η_t adaptively such that η_t becomes larger when the current target \mathbf{w}^* is far from \mathbf{w}_t and becomes small when it is close to \mathbf{w}_t adaptively. However, no definite scheme was analyzed there. Sompolinsky et al. [1995] proposed an excellent scheme of an adaptive choice of η_t for a deterministic dichotomy neural networks. We extend their idea to be applicable to stochastic cases, where the Riemannian structure plays a role.

We assume that $l(\mathbf{x}, \mathbf{z}; \mathbf{w})$ is differentiable with respect to \mathbf{w} . (The non-differentiable case is usually more difficult to analyze. Sompolinsky et al [1995] treated this case.) We moreover treat the realizable teacher so that $L(\mathbf{w}^*) = 0$.

We propose the following learning scheme:

$$\hat{\mathbf{w}}_{t+1} = \hat{\mathbf{w}}_t - \eta_t \tilde{\nabla} l(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}; \hat{\mathbf{w}}_t) \quad (14)$$

$$\eta_{t+1} = \eta_t + \alpha \eta_t [\beta l(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}; \hat{\mathbf{w}}_t) - \eta_t], \quad (15)$$

where α and β are constants. We try to analyze the dynamical behavior of learning by using the continuous version of the algorithm,

$$\frac{d}{dt} \hat{\mathbf{w}}_t = -\eta_t \tilde{\nabla} l(\mathbf{x}_t, \mathbf{z}_t; \hat{\mathbf{w}}_t), \quad (16)$$

$$\frac{d}{dt} \eta_t = \alpha \eta_t [\beta l(\mathbf{x}_t, \mathbf{z}_t; \hat{\mathbf{w}}_t) - \eta_t]. \quad (17)$$

In order to show the dynamical behavior of $(\hat{\mathbf{w}}_t, \eta_t)$, we use the averaged version of the above equation with respect to the current input-output pair $(\mathbf{x}_t, \mathbf{z}_t)$. We introduce the squared error variable

$$e_t = \frac{1}{2} (\mathbf{w}_t - \mathbf{w}^*)' G^* (\mathbf{w}_t - \mathbf{w}^*). \quad (18)$$

By using the average and continuous time version

$$\dot{\mathbf{w}}_t = -\eta_t G^{-1}(\mathbf{w}_t) \left\langle \frac{\partial}{\partial \mathbf{w}} l(\mathbf{x}_t, \mathbf{z}_t; \mathbf{w}_t) \right\rangle,$$

$$\dot{\eta}_t = \alpha \eta_t \{ \beta \langle l(\mathbf{x}_t, \mathbf{z}_t; \mathbf{w}_t) \rangle - \eta_t \},$$

where $\dot{\cdot}$ denotes d/dt and $\langle \cdot \rangle$ the average over the current (\mathbf{x}, \mathbf{z}) , we have

$$\dot{e}_t = -2\eta_t e_t, \quad (19)$$

$$\dot{\eta}_t = \alpha \beta \eta_t e_t - \alpha \eta_t^2. \quad (20)$$

The behavior of the above equation is interesting : The origin $(0, 0)$ is its attractor. However, the basin of attraction has a fractal boundary. Anyway, starting from an adequate initial value, it has the solution of the form

$$e_t \approx \frac{1}{\beta} \left(\frac{1}{2} - \frac{1}{\alpha} \right) \frac{1}{t}, \quad (21)$$

$$\eta_t \approx \frac{1}{2t}. \quad (22)$$

This proves the $1/t$ convergence rate of the generalization error, that is optimal in order for any estimator $\hat{\mathbf{w}}_t$ converging to \mathbf{w}^* .

6 Riemannian structures of simple perceptrons

We first study the parameter space S of simple perceptrons to obtain an explicit form of the metric G and its inverse G^{-1} . This suggests how to calculate the metric in the parameter space of multilayer perceptrons.

Let us consider a simple perceptron with input \mathbf{x} and output z . Let \mathbf{w} be its connection weight vector. For the analog stochastic perceptron, its input-output behavior is described by $z = f(\mathbf{w}'\mathbf{x}) + n$, where n denotes a random noise subject to the normal distribution $N(0, \sigma^2)$ and f is the hyperbolic tangent,

$$f(u) = \frac{1 - e^{-u}}{1 + e^{-u}}.$$

In order to calculate the metric G explicitly, let $\mathbf{e}_{\mathbf{w}}$ be the unit column vector in the direction of \mathbf{w} in the Euclidean space \mathbf{R}^n ,

$$\mathbf{e}_{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|},$$

where $\|\mathbf{w}\|$ is the Euclidean norm. We then have the following theorem.

Theorem 2. The Fisher metric G and its inverse G^{-1} are given by

$$G(\mathbf{w}) = c_1(w)I + \{c_2(w) - c_1(w)\}\mathbf{e}_{\mathbf{w}}\mathbf{e}'_{\mathbf{w}}, \quad (23)$$

$$G^{-1}(\mathbf{w}) = \frac{1}{c_1(w)}I + \left(\frac{1}{c_2(w)} - \frac{1}{c_1(w)} \right) \mathbf{e}_{\mathbf{w}}\mathbf{e}'_{\mathbf{w}}. \quad (24)$$

where $w = \|\mathbf{w}\|$ (Euclidean norm) and $c_1(w)$ and $c_2(w)$ are given by

$$c_1(w) = \frac{1}{4\sqrt{2\pi}\sigma^2} \int \{f^2(w\varepsilon) - 1\}^2 \exp\left\{-\frac{1}{2}\varepsilon^2\right\} d\varepsilon, \quad (25)$$

$$c_2(w) = \frac{1}{4\sqrt{2\pi}\sigma^2} \int \{f^2(w\varepsilon) - 1\}^2 \varepsilon^2 \exp\left\{-\frac{1}{2}\varepsilon^2\right\} d\varepsilon. \quad (26)$$

Theorem 3. The Jeffrey prior is given by

$$\sqrt{|G(\mathbf{w})|} = \frac{1}{V_n} \sqrt{c_2(w)\{c_1(w)\}^{n-1}}. \quad (27)$$

7 The natural gradient for blind separation of mixtured signals

Let $\mathbf{s} = (s_1, \dots, s_n)$ be n source signals which are n independent random variables. We assume that their n mixtures $\mathbf{x} = (x_1, \dots, x_n)$,

$$\mathbf{x} = A\mathbf{s} \quad (28)$$

are observed. Here, A is a matrix. When \mathbf{s} is time serieses, we observe $\mathbf{x}(1), \dots, \mathbf{x}(t)$. The problem of blind separation is to estimate $W = A^{-1}$ adaptively from $\mathbf{x}(t)$, $t = 1, 2, 3, \dots$ without knowing $\mathbf{s}(t)$ nor A . We can then recover original \mathbf{s} by

$$\mathbf{y} = \hat{W}\mathbf{x} \quad (29)$$

when $\hat{W} = A^{-1}$. Let $W \in Gl(n)$, that is a nonsingular $n \times n$ -matrix, and $\phi(W)$ be a scalar function. This is given by a measure of independence such as $\phi(W) = KL[\hat{p}(\mathbf{y}); p(\mathbf{y})]$, which is represented by the expectation of a loss function. We define the natural gradient of $\phi(W)$.

Now we return to our manifold $Gl(n)$ of matrices. It has the Lie group structure : Any $A \in Gl(n)$ maps $Gl(n)$ to $Gl(n)$ by $W \rightarrow WA$. We impose that the Riemannian structure should be invariant by this operation A .

We can then prove that the natural gradient in this case is

$$\tilde{\nabla}\phi = \nabla\phi W'W. \quad (30)$$

The natural gradient works surprisingly well for adaptive blind signal separation Amari et al. [1995], Cardoso and Laheld [1996].

References

- [1] S. Amari. Theory of adaptive pattern classifiers, *IEEE Trans.*, **EC-16**, No.3, 299–307, 1967.
- [2] S. Amari. *Differential-Geometrical Methods in Statistics, Lecture Notes in Statistics*, vol.28, Springer, 1985.
- [3] S. Amari. Information geometry of the EM and em algorithms for neural networks, *Neural Networks*, **8**, No.9, 1379–1408, 1995.
- [4] S. Amari, A. Cichocki and H.H. Yang. A new learning algorithm for blind signal separation, in *NIPS'95*, vol.8, 1996, MIT Press, Cambridge, Mass.
- [5] S. Amari, K. Kurata, H. Nagaoka. Information geometry of Boltzmann machines, *IEEE Trans. on Neural Networks*, **3**, 260–271, 1992.
- [6] J. F. Cardoso and Beate Laheld. Equivariant adaptive source separation, to appear *IEEE Trans. on Signal Processing*, 1996.
- [7] T. M. Heskes and B. Kappen. Learning processes in neural networks, *Physical Review A*, **440**, 2718–2726, 1991.
- [8] D. Rumelhart, G.E. Hinton and R. J. Williams. Learning internal representation, in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, **1**, *Foundations*, MIT Press, Cambridge, MA, 1986.
- [9] H. Sompolinsky, N. Barkai and H. S. Seung. On-line learning of dichotomies: algorithms and learning curves, *Neural Networks: The statistical Mechanics Perspective*, Proceedings of the CTP-PBSRI Joint Workshop on Theoretical Physics, J.-H. Oh et al eds, 105–130, 1995.