

Gradient Learning in A Non-Euclidean Space

Tsung-Wei Chiang (Jimmy)

Natinal Taiwan University

d98942029@ntu.edu.tw

September 11, 2018

Motivation:

Parameter Space from Euclidean to Non-Euclidean

- Previous Assumption: Learning takes places in a parameter space which is Euclidean.
- New Assumption: Learning takes places in a parameter space which is Non-Euclidean.
- What if the parameter space is Non-Euclidean $[1, 2, 3]$, e.g., Reimannian?
- (Ordinary) Gradient Descent Methods (including SGD) in deep learning
- Natural Gradient Descent Methods

- Example: a regression problem,
- input signal: \mathbf{x} , generate randomly
- desired response: $f(\mathbf{x})$
- teacher signal: y
- random noise: $\epsilon \sim p_\epsilon(\epsilon)$, e.g., Gaussian
- noisy version of the desired output:

$$y = f(\mathbf{x}) + \epsilon \quad (1)$$

- unknown joint probability distribution: $p(\mathbf{x})$
- The task of a learning machine is to estimate the desired output mapping $f(\mathbf{x})$ by using the available examples of input-output pairs $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, T\}$ (training examples)

- $f(\mathbf{x}, \xi)$: a parametrized family of functions as candidates for the desired output
- ξ : a vector parameter. The set of ξ is a parameter space.
- We search for the optimal $\hat{\xi}$ that approximates the ture $f(\xi)$ by using training examples D .
- Regression Problem: when y takes an analog value.
- Pattern Recognition: when y is decrete, say, binary.

Overview

- The instantaneous loss function (of processing \mathbf{x} by machine $f(\mathbf{x}, \xi)$), e.g.:

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

- Generalization error, the loss function function (of machine ξ is the expectation of the instantaneous loss over all possible pairs (\mathbf{x}, y) :

$$L(\xi) = E_{p(\mathbf{x}, y)} \{l(\mathbf{x}, y, \xi)\} \quad (3)$$

where the expectation is taken with respect to the unknown joint probability distribution $p(\mathbf{x}, y)$.

- Training error:

$$L_{\text{train}}(\xi) = \frac{1}{T} \sum_{t=1}^T l(\mathbf{x}_t, y_t; \xi_t) \quad (4)$$

since we do not know $p(\mathbf{x}, y)$, we use the average over the training data.

- Since we do not know L , we minimize the training error L_{train} to obtain $\hat{\xi}$.
- A regularization term may be added to L_{train} .

- On-line Learning: Modifying the current candidate ξ_t at time t to obtain ξ_{t+1} at the next time $t + 1$ based on the current training example (\mathbf{x}_t, y_t) so as to decrease the instantaneous loss $l(\mathbf{x}_t, y_t, \xi_t)$.
- Usually, the negative of the gradient is used to update ξ_t :

$$\xi_{t+1} = \xi_t - \eta_t \nabla_{\xi} l(\mathbf{x}_t, y_t, \xi_t) \quad (5)$$

- ∇ : the gradient with respect to ξ
- η_t : learning rate, which may depend on t .

On-line Learning

- Since training data (\mathbf{x}_t, y_t) are given one by one, the change is a random variable,

$$\Delta \xi_t = -\eta_t \nabla_{\xi} l(\mathbf{x}_t, y_t, \xi_t) \quad (6)$$

- The expectation of $\nabla_{\xi} l(\mathbf{x}_t, y_t, \xi_t)$ equals $\nabla_{\xi} L(\xi_t)$.
- Stochastic Gradient Descent (SGD) Learning:

$$\Delta \xi_t = -\eta_t \nabla_{\xi} l(\mathbf{x}_t, y_t, \xi_t) \quad (7)$$

- Gradient Descent (GD) Learning:

$$E\{\Delta \xi_t\} = -\eta_t E\{\nabla_{\xi} l(\mathbf{x}_t, y_t, \xi_t)\} = -\eta_t \nabla_{\xi} L(\xi_t) \quad (8)$$

- The change of $\Delta \xi_t$ is random but its expectation is in the direction of $-\nabla_{\xi} L(\xi_t)$.
- Well established as the back-propagation learning method.

Gradient descent of expected loss L v.s. stochastic gradient descent of l .

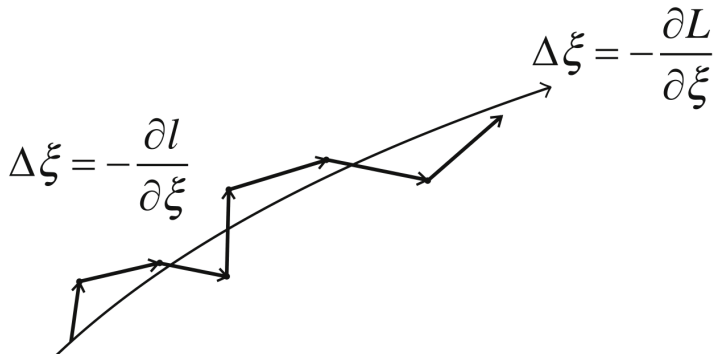


Figure: See [4].

- Batch Learning: an iterative method which uses all the training data for modifying ξ_t at one step, such that ξ_t is modified to ξ_{t+1} by
- The update

$$\xi_{t+1} = \xi_t - \eta_t \frac{1}{T} \sum_{i=1}^T \nabla_{\xi} l(\mathbf{x}_i, y_i, \xi_t) \quad (9)$$

What is Natural Gradient

- It is believed that the gradient

$$\nabla_{\xi} L(\xi) = \frac{\partial L(\xi)}{\partial \xi} \quad (10)$$

is the direction of the steepest change of $L(\xi)$.

- In a geographical map with contour lines, the steepest direction is given by the gradient of the height function $H(\xi)$, i.e. $\nabla_{\xi} H(\xi)$ is orthogonal to contour lines.
- However, this is true only when an orthonormal coordinate system is used in a Euclidean space.
- Ordinary Gradient: steepest descent direction in Euclidean Manifolds
- Natural Gradient: steepest descent direction in Riemannian Manifolds

What is Natural Gradient

- It is believed that the gradient

$$\nabla_{\xi} L(\xi) = \frac{\partial L(\xi)}{\partial \xi} \quad (11)$$

is the direction of the steepest change of $L(\xi)$.

- In a geographical map with contour lines, the steepest direction is given by the gradient of the height function $H(\xi)$, i.e. $\nabla_{\xi} H(\xi)$ is orthogonal to contour lines.
- However, this is true only when an orthonormal coordinate system is used in a Euclidean space.
- Ordinary Gradient: steepest descent direction in Euclidean Manifolds
- Natural Gradient: steepest descent direction in Riemannian Manifolds

A Riemannian Manifold

- A Riemannian manifold can be intuitively seen as one kind of high dimensional ($\geq 4D$, or k -manifold, $k \geq 3$) surface, a hypersurface.

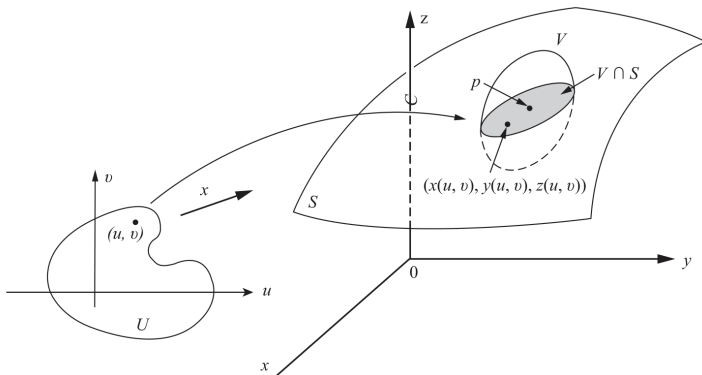


Figure: A 2-manifold, 3D surface. See [5].

A Riemannian Manifold

- In a Riemannian manifold, the square of local distance between two nearby points ξ and $\xi + d\xi$ is given by the quadratic form:

$$ds^2 = g_{ij} d\xi^i d\xi^j \quad (12)$$

- Einstein convention is applied here.
- $\mathbf{G} = [g_{ij}]$: a Riemannian metric tensor.
- Let us change the current point ξ to $\xi + d\xi$, and see how the value of $L(\xi)$ changes, depending on the direction $d\xi$. We search for the direction in which L changes most rapidly.

Natural Gradient

- In order to make a fair comparison, the step-size of $d\xi$ should have the same magnitude in all directions, so that the length of $d\xi$ should be the same,

$$g_{ij}(\xi)d\xi^i d\xi^j = \epsilon^2 \quad (13)$$

- ϵ : a small constant
- We put $d\xi = \epsilon \mathbf{a}$ and require that

$$|\mathbf{a}|^2 = g_{ij}a^i a^j = 1 \quad (14)$$

- The Steepest direction of L is the maximizer of

$$L(\xi + d\xi) - L(\xi) = \epsilon \nabla L(\xi) \cdot \mathbf{a} \quad (15)$$

under the constraint (14)

Natural Gradient

- Using the variational method, we obtain:

$$\underset{\mathbf{a}}{\text{maximize}} \nabla_{\xi} L(\xi) \cdot \mathbf{a} - \lambda g_{ij} a^i a^j \quad (16)$$

- This is a quadratic problem and the steepest direction is obtained as

$$\mathbf{a} \propto \mathbf{G}^{-1} \nabla_{\xi} L(\xi) \quad (17)$$

- The natural gradient or Riemannian gradient of L is

$$\widetilde{\nabla}_{\xi} L(\xi) = \mathbf{G}^{-1}(\xi) \nabla_{\xi} L(\xi) \quad (18)$$

- The natural gradient operator is

$$\widetilde{\nabla}_{\xi} = \mathbf{G}^{-1}(\xi) \nabla_{\xi} \quad (19)$$

Our goal:

$$\underset{\mathbf{a}}{\text{maximize}} \nabla_{\xi} L(\xi) \cdot \mathbf{a} - \lambda g_{ij} a^i a^j \quad (20)$$

By the Lagrangian method, we take

$$\frac{\partial}{\partial a_i} \{ \nabla_{\xi} L(\xi)^{\top} \mathbf{a} - \lambda \mathbf{a}^{\top} \mathbf{G} \mathbf{a} \} = 0 \quad (21)$$

Then, we get

$$\nabla_{\xi} L(\xi) = 2\lambda \mathbf{G} \mathbf{a} \quad (22)$$

$$\mathbf{a} = \frac{1}{2\lambda} \mathbf{G}^{-1} \nabla_{\xi} L(\xi) \quad (23)$$

Using the constraint $|\mathbf{a}|^2 = g_{ij} a^i a^j = 1$, λ is determined.

Natural Gradient: Discussions

- Special Case: when an orthonormal coordinate system is used in a Euclidean space, we have

$$g_{ij}(\boldsymbol{\xi}) = \delta_{ij} \quad (24)$$

- On-line learning mode, the natural gradient decent method:

$$\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t - \eta_t \widetilde{\nabla}_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_t) \quad (25)$$

- Batch learning mode,

$$\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t - \eta_t \frac{1}{T} \sum_{i=1}^T \widetilde{\nabla}_{\boldsymbol{\xi}} l(\mathbf{x}_i, y_i, \boldsymbol{\xi}_t) \quad (26)$$

Natural Gradient: Discussions

In the case of statistical estimation,

- The loss function and the Riemannian metric \mathbf{G} are related.
- The Fisher information is a Riemannian metric

$$g_{ij} = E \left[\frac{\partial \log p(\mathbf{x}, \xi)}{\partial \xi_i} \frac{\partial \log p(\mathbf{x}, \xi)}{\partial \xi_j} \right] \quad (27)$$

- The loss function L uses the same log likelihood $\log p(\mathbf{x}, \xi)$.
- The Riemannian \mathbf{G} uses the same log likelihood $\log p(\mathbf{x}, \xi)$.
- The natural gradient method is a version of the Gauss-Newton method.

Natural Gradient: Discussions

In the case of independent component analysis,

- The loss function and the Riemannian metric \mathbf{G} are NOT related.
- The natural gradient learning method is useful in such case, too.
- Parameter space is a set of mixing matrices.
- The loss function L is measured by the degree of independence of unmixed signals.
- The Riemannian \mathbf{G} is measured by the invariant metric of Lie group.

In the case of independent component analysis,

- Deep learning: [Roux et al. 2007; Ollivier 2015]
- Reinforcement learning: as a policy natural gradient [Kakade 2002; Peters and Schaal 2008; Morimura et al. 2009]
- Optimization: stochastic relaxation technique [Malagò and Pistone 2014; Malagò et al. 2013; Yi et al. 2009; see also Hansen and Ostermeier 2001]

Natural Gradient: Property 1

Natural gradient learning is Fisher efficient.

Theorem

The estimator obtained by on-line natural gradient learning

$$\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t - \eta_t \widetilde{\nabla}_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_t) \quad (28)$$

is Fisher efficient, attaining the Cramér-Rao lower bound asymptotically.

Natural Gradient: Property 1, Proof

- Let the error covariance matrix of the estimator at time t be

$$\mathbf{V}_{t+1} = E[(\boldsymbol{\xi}_{t+1} - \boldsymbol{\xi}_0)(\boldsymbol{\xi}_{t+1} - \boldsymbol{\xi}_0)^T] \quad (29)$$

- $\boldsymbol{\xi}_0$ is the true value of $\boldsymbol{\xi}$.
- We expand the loss at $\boldsymbol{\xi}_t$ as

$$\nabla_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_t) = \nabla_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_0) + \nabla_{\boldsymbol{\xi}} \nabla_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_0) \cdot (\boldsymbol{\xi}_t - \boldsymbol{\xi}_0) \quad (30)$$

- Subtracting $\boldsymbol{\xi}_0$ from both sides of (28) and substituting it in (29), we have

$$\mathbf{V}_{t+1} = \mathbf{V}_t - \frac{2}{t} \mathbf{V}_t + \frac{1}{t^2} \mathbf{G}^{-1} + O\left(\frac{1}{t^3}\right) \quad (31)$$

where

$$E[\nabla_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_0)] = 0 \quad (32)$$

$$E[\nabla_{\boldsymbol{\xi}} \nabla_{\boldsymbol{\xi}} l(\mathbf{x}_t, y_t, \boldsymbol{\xi}_0)] = \mathbf{G}(\boldsymbol{\xi}_0) \quad (33)$$

- Note that

$$\mathbf{G}(\boldsymbol{\xi}_t) = \mathbf{G}(\boldsymbol{\xi}_0) + O\left(\frac{1}{t}\right) \quad (34)$$

- Then the solution of (31) is asymptotically

$$\mathbf{V}_t = \frac{1}{t} \mathbf{G}^{-1} \quad (35)$$

which prove the theorem.

Natural Gradient: Property 2

- Consider a regression problem, the output is written as

$$y = f(\mathbf{x}, \boldsymbol{\xi}) + \epsilon \quad (36)$$

- First, we consider a simple perceptron, where f is written as

$$f(\mathbf{x}, \boldsymbol{\xi}) = \phi(\mathbf{w} \cdot \mathbf{x}) \quad (37)$$

- Here, we neglect the bias term for simplicity.
- The parameter is a vector $\boldsymbol{\xi} = \mathbf{w}$ and the activation function ϕ is a sigmoid function,

$$\phi(u) = \tanh(u) \quad (38)$$

- The gradient is

$$\nabla l(\mathbf{x}, y, \mathbf{w}) = -(y - f)\phi'(\mathbf{w} \cdot \mathbf{x})\mathbf{x} \quad (39)$$

- When the absolute value of \mathbf{w} is large, $\phi(\mathbf{w} \cdot \mathbf{x})$ saturates for most \mathbf{x} , becoming nearly equal $+1$ or -1 .
- This is the saturation problem, where the gradient almost equal to 0 because $\phi' \approx 0$, and the ordinary SGD learning becomes slow.
- This is not serious in the case of a simple perceptron, but is serious in the case of multilayer perceptrons used in deep learning, where $f(\mathbf{x}, \boldsymbol{\xi})$ is composed of a concatenation of many f . In MLP, We may write the output as

$$f(\mathbf{x}, \boldsymbol{\xi}) = \phi(\mathbf{W}_k \phi(\mathbf{W}_{k-1} \phi(\cdots \phi(\mathbf{W}_1 \mathbf{x}))) \quad (40)$$

where $\boldsymbol{\xi} = (\mathbf{W}_1, \cdots, \mathbf{W}_k)$

- Its derivative with respect to \mathbf{W}_1 , for example, includes the product of many ϕ' . Hence, it is almost vanishing in many cases. This is considered as a flaw of back-propagation in deep learning.

Natural Gradient: Property 2

- The natural gradient learning method is free of such a saturation problem. The gradient is written as

$$\nabla l(\mathbf{x}, y, \xi) = -(y - f) \nabla f(\mathbf{x}, \xi) \quad (41)$$

- The magnitude of the ordinary gradient would be very small in many cases but the natural gradient is different.

Natural Gradient: Property 2

Natural gradient is Saturation Free.

Theorem

The magnitude of the natural gradient is given by

$$E \left[\left\| \tilde{\nabla} l \right\|^2 \right] = \text{tr} \left(\overline{\mathbf{G}}(\xi) \overline{\mathbf{G}}^{-1}(\xi) \right) \quad (42)$$

where

$$\overline{\mathbf{G}}(\xi) = E_{p(\mathbf{x}, y, \xi_0)} \left[\nabla_{\xi} l(\mathbf{x}, \xi) l(\mathbf{x}, \xi)^{\top} \right] \quad (43)$$

It dose not vanish even when ϕ' is small. Moreover,

$$E \left[\left\| \tilde{\nabla} l \right\|^2 \right] \approx k \quad (44)$$

in a neighborhood of the optimal ξ_0 , where k is the dimension of ξ .

Natural Gradient: Property 2, Proof

Firstly, from

$$\nabla_{\xi} l(\mathbf{x}, \xi) = \mathbf{G}^{-1}(\mathbf{x}, \xi) \nabla_{\xi} l(\mathbf{x}, \xi) \quad (45)$$

We have

$$E \left[\left\| \tilde{\nabla} l \right\|^2 \right] = E_{p(\mathbf{x}, y, \xi_0)} \left[\text{tr} \mathbf{G}(\xi) \mathbf{G}^{-1}(\xi) \nabla_{\xi} l(\mathbf{x}, \xi) l(\mathbf{x}, \xi)^{\top} \mathbf{G}^{-1}(\xi) \right] \quad (46)$$

which completes the proof.

Secondly, when $\xi = \xi_0$, we easily have (44).

References I



S.-i. Amari, "Natural gradient works efficiently in learning," *Neural Computation*, vol. 10, no. 2, pp. 251–276, 1998. [Online]. Available: <https://doi.org/10.1162/089976698300017746>



S. Fiori, "Extended hamiltonian learning on riemannian manifolds: Theoretical aspects," *IEEE Transactions on Neural Networks*, vol. 22, no. 5, pp. 687–700, May 2011.



S. Fiori, "Extended hamiltonian learning on riemannian manifolds: Numerical aspects," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 23, no. 1, pp. 7–21, Jan 2012.



S.-i. Amari, *Information Geometry and Its Applications*, 1st ed. Springer Publishing Company, Incorporated, 2016.



M. do Carmo, *Differential Geometry of Curves and Surfaces*. Prentice-Hall, 1976. [Online]. Available: <https://books.google.com.tw/books?id=1v0YAQAAIAAJ>

The End