Contents lists available at ScienceDirect

Neural Networks

journal homepage: www.elsevier.com/locate/neunet

Non-monotonic convergence of online learning algorithms for perceptrons with noisy teacher

Kazushi Ikeda^{a,*}, Arata Honda^{a,1}, Hiroaki Hanzawa^{a,2}, Seiji Miyoshi^b

^a Nara Institute of Science and Technology, Ikoma, Nara, Japan

^b Kansai University, Suita, Osaka, Japan

ARTICLE INFO

Article history: Received 22 March 2017 Received in revised form 29 November 2017 Accepted 9 February 2018 Available online 21 February 2018

Keywords: Learning curve Perceptron Online learning Statistical mechanics Asymptotic analysis

1. Introduction

Statistical mechanical methods can apply to problems in information science such as neural networks (Nishimori, 2001), communication theory (Tanaka, 2002), and adaptive filters (Miyoshi & Kajikawa, 2013). One successful application is the analyses of the perceptron learning algorithm (Biehl & Schwarze, 1992; Rosenblatt, 1961) and its variations (Hara & Okada, 2004; Inoue & Nishimori, 1997; Miyoshi, Hara, & Okada, 2005; Miyoshi & Okada, 2006a, b; Uezu, Miyoshi, Izuo, & Okada, 2007).

The perceptron learning is an online learning algorithm where the student updates its weight vector of a linear dichotomy according to the teacher's signal (Rosenblatt, 1961). Biehl and Schwarze (1992) introduced the statistical mechanics to the analysis of the perceptron learning and Inoue and Nishimori (1997) applied the method to the AdaTron learning in unlearnable cases. Hara and Okada (2004) discussed the perceptron learning with a margin and Miyoshi and his colleagues extended the analysis to the ensemble learning and/or noisy cases (Miyoshi et al., 2005; Miyoshi & Okada, 2006a, b; Uezu et al., 2007).

In this paper, we consider the case where the teacher has noise in its output while the student does not. In this case, the learning

² H. Hanzawa is currently with Yahoo! Japan.

https://doi.org/10.1016/j.neunet.2018.02.009 0893-6080/© 2018 Elsevier Ltd. All rights reserved.

ABSTRACT

Learning curves of simple perceptron were derived here. The learning curve of the perceptron learning with noisy teacher was shown to be non-monotonic, which has never appeared even though the learning curves have been analyzed for half a century. In this paper, we showed how this phenomenon occurs by analyzing the asymptotic property of the perceptron learning using a method in systems science, that is, calculating the eigenvalues of the system matrix and the corresponding eigenvectors. We also analyzed the AdaTron learning and the Hebbian learning in the same way and found that the learning curve of the AdaTron learning is non-monotonic whereas that of the Hebbian learning is monotonic.

© 2018 Elsevier Ltd. All rights reserved.

curve, which is defined as the average prediction error, is not monotonically decreasing but has an overshoot, differently from other cases analyzed so far (Ikeda, Hanzawa, & Miyoshi, 2013). Although an analysis for this problem was partially given by some of the authors (Ikeda et al., 2013), some part was given not theoretically but numerically.

This paper gives a theoretically rigorous and complete analvsis above. In addition, we extend the analysis to other online algorithms for perceptrons, that is, the AdaTron learning and the Hebbian learning. Our analysis consisted of three steps. In the first step, we applied the statistical mechanical method to our problem, i.e., we introduced three order parameters assuming the thermodynamic limit, and derived a system of differential equations for the three algorithms. In the second step, we calculated the ensemble averages that appeared in the differential equations for each algorithm using Gaussian approximations. Note this had not been derived analytically yet in Ikeda et al. (2013). In the last step, we applied an asymptotic analysis to our dynamical system, i.e., we linearized the equations around their convergence point and analyzed their behaviors by the eigenvalues and eigenvectors of the state-transition matrix a.k.a. the system matrix. The three steps elucidated how and why the overshoot phenomenon occurs.

The remainder of this paper is organized as follows. Section 2 formulates the problem we treated. Sections 3–5 are devoted to the three steps, that is, statistical mechanical analysis, the calculation of the ensemble averages and the asymptotic analysis of the system, respectively. We conclude the paper in Section 6.







^{*} Corresponding author.

E-mail addresses: kazushi@is.naist.jp (K. Ikeda), arata.honda@excite.jp (A. Honda), h.hanzawax68@gmail.com (H. Hanzawa), miyoshi@kansai-u.ac.jp

⁽S. Miyoshi).

¹ A. Honda is currently with Excite.

2. Problem statement

Two linear perceptrons are treated: a teacher and a student, whose connection weights are $B = (B_1, \ldots, B_N) \in \mathbb{R}^N$ and J = $(J_1, \ldots, J_N) \in \mathbb{R}^N$, respectively. The initial value of each of the components is independently drawn from the normal distribution *N*(0, 1), that is,

$$\langle B_i \rangle = 0, \qquad \langle (B_i)^2 \rangle = 1, \qquad (1)$$

$$\langle I_i \rangle = 0, \qquad \langle (I_i)^2 \rangle = 1, \qquad (2)$$

where $\langle \cdot \rangle$ denotes the mean of \cdot , as was in Nishimori (2001).

The *m*th input vector $x^m = (x_1^m, \ldots, x_N^m) \in \mathbb{R}^N$ is independently drawn from the N-dimensional normal distribution N(0, I/N) and the corresponding output y^m of the teacher is produced as

$$y^m = \operatorname{sgn}(v_m), \qquad v_m = B \cdot x^m + n_B^m, \qquad (3)$$

where n_B^m is an observation noise obeying $N(0, \sigma_B^2)$. The learning rule is either the standard perceptron learning (Biehl & Schwarze, 1992; Nishimori, 2001; Rosenblatt, 1961), the AdaTron learning (Nishimori, 2001), or the Hebbian learning (Nishimori, 2001). In the perceptron learning, given the *m*th input vector x^m , the student updates its weight vector J^m as

$$J^{m+1} = J^m + f^m x^m, (4)$$

$$f^{m} = \eta y^{m} \Theta(-y^{m} J^{m} \cdot x^{m}),$$
(5)

where η is a learning coefficient and $\Theta(\cdot)$ is the Heaviside function,

$$\Theta(t) = \begin{cases} 1 & t \ge 0, \\ 0 & t < 0. \end{cases}$$
(6)

This means that it updates its weight vector when its output does not coincide with the teacher's one.

In the AdaTron learning and the Hebbian learning, the update functions f^m are changed to

$$f^{m} = \eta y^{m} J^{m} x^{m} \Theta(-y^{m} J^{m} \cdot x^{m}),$$

$$f^{m} = \eta y^{m},$$
(8)

respectively.

As the learning proceeds and m increases, the weight vector, I^m , of the student approaches the teacher's one, B. The problem of learning curves is to evaluate how fast the covariance coefficient between J^m and B,

$$R^m = \frac{B \cdot J^m}{\|B\| \|J^m\|},\tag{9}$$

approaches unity in noiseless cases and another value in noisy cases (0.70 in Fig. 1, for example).

3. Statistical mechanical analysis

3.1. Theory

The method to derive the learning curve of the student is essentially the same as Nishimori (2001). We introduce auxiliary order parameters, R^m in (9) and

$$l^m = \|J^m\|/\sqrt{N},\tag{10}$$

and consider the thermodynamic limit, $N, m \rightarrow \infty$ and m/N = t. Then,

$$||B|| = \sqrt{N},$$
 $||J^0|| = \sqrt{N},$ $||x^m|| = 1,$ (11)

hold and the random vector of the inner products, (u, v), where

$$v^m = B \cdot x^m, \qquad u^m l^m = J^m \cdot x^m \tag{12}$$



Fig. 1. Dynamics of *R*. $\sigma_B^2 = 1.0$, $N = 10^4$, $\eta = 0.1, ..., 2.0$, plots: experiments, lines: theory (modified from Ikeda et al., 2013).

obeys the two-dimensional normal distribution $N(0, \Sigma)$ where

$$\Sigma = \begin{pmatrix} 1 & R^m \\ R^m & 1 \end{pmatrix}.$$
 (13)

By self-averaging and omitting the step index m in (4) hereafter, we get the simultaneous differential equations of the order parameters.

$$\dot{l} = \langle f u \rangle + \frac{\langle f^2 \rangle}{2l},\tag{14}$$

$$\dot{R} = \frac{\langle f v \rangle - \langle f u \rangle R}{l} - \frac{R}{2l^2} \langle f^2 \rangle, \qquad (15)$$

where $\langle \cdot \rangle$ expresses the average over (u, v) and $n_B \sim N(0, \sigma_B^2)$ (Nishimori, 2001).

3.2. Experiments

To confirm the validity of the theory above, we conducted some computer simulations of the perceptron learning under the condition in Section 2. The experimental values of *R* coincided well with the theoretical values for any learning coefficient, η (Fig. 1).

As a result of the experiments, the value of *R* converged to 0.70 for any η due to the noise on the teacher's output. One notable property was that R was not monotonically increasing but had an overshoot. This overshoot phenomenon does not occur in other cases analyzed so far (Hara & Okada, 2004; Inoue & Nishimori, 1997; Miyoshi et al., 2005; Miyoshi & Okada, 2006a; Nishimori, 2001).

A quantitative analysis of this phenomenon is given using an asymptotic dynamical system theory in Section 5.

4. Calculation of ensemble averages

The ensemble averages $\langle f v \rangle$, $\langle f u \rangle$ and $\langle f^2 \rangle$ in (14) and (15) are difficult to calculate analytically, in general. In fact, we calculated those for the perceptron learning numerically (Ikeda et al., 2013). However, we theoretically derived the ensemble averages for the perceptron learning. In addition, we also calculated those for the AdaTron learning and the Hebbian learning, which will be given below

The ensemble averages $\langle f v \rangle$, $\langle f u \rangle$ and $\langle f^2 \rangle$ for the perceptron learning are expressed as

$$\langle fv \rangle = \langle \eta \Theta(-u(v+n_B)) \operatorname{sgn}(v+n_B)v \rangle$$

= $\eta \int_{-\infty}^{\infty} dn_B \int_{-n_B}^{\infty} dv \int_{-\infty}^{0} du P(u,v) P(n_B)v$

Download English Version:

https://daneshyari.com/en/article/6862999

Download Persian Version:

https://daneshyari.com/article/6862999

Daneshyari.com