

## Center Manifold Analysis of Plateau Phenomena Caused by Degeneration of Three-Layer Perceptron

Daiji Tsutsui

*d-tsutsui@cr.math.sci.osaka-u.ac.jp*

*Osaka University, Toyonaka-shi, Osaka 560-0043, Japan*

A hierarchical neural network usually has many singular regions in the parameter space due to the degeneration of hidden units. Here, we focus on a three-layer perceptron, which has one-dimensional singular regions comprising both attractive and repulsive parts. Such a singular region is often called a Milnor-like attractor. It is empirically known that in the vicinity of a Milnor-like attractor, several parameters converge much faster than the rest and that the dynamics can be reduced to smaller-dimensional ones. Here we give a rigorous proof for this phenomenon based on a center manifold theory. As an application, we analyze the reduced dynamics near the Milnor-like attractor and study the stochastic effects of the online learning.

### 1 Introduction ---

A three-layer perceptron is one of the simplest hierarchical learning machines. Mathematically, it is defined by the function

$$f_{(d)}(x; \theta) = \sum_{i=1}^d v_i \varphi(\mathbf{w}_i \cdot \mathbf{x} + b_i), \quad \mathbf{x} \in \mathbb{R}^n, \quad (1.1)$$

$$\theta = (\mathbf{w}_1, \dots, \mathbf{w}_d, b_1, \dots, b_d, \mathbf{v}_1, \dots, \mathbf{v}_d),$$

where  $\theta$  is a system parameter with  $\mathbf{w}_1, \dots, \mathbf{w}_d \in \mathbb{R}^n$  being the weight vectors for the second layer,  $b_1, \dots, b_d \in \mathbb{R}$  the bias terms for the second,  $\mathbf{v}_1, \dots, \mathbf{v}_d \in \mathbb{R}^m$  the weight vectors for the third, and  $\varphi$  an activation function. Throughout this article, we assume that the activation function is twice differentiable. Figure 1 is a schematic diagram of the three-layer perceptron. We shall call the function 1.1 an  $(n-d-m)$ -perceptron. The numbers  $n$  and  $m$  are fixed at the outset as the sizes of input and output vectors, while the number  $d$  of hidden units can be varied in our analysis. For notational simplicity, we incorporate the bias  $b$  in the weight  $\mathbf{w}$  as  $\mathbf{w} = (b, w^1, \dots, w^n)$ , and accordingly, we enlarge  $\mathbf{x}$  as  $\mathbf{x} = (1, x_1, \dots, x_n)$ . By using these conventions, we obtain the abridged presentation of the three-layer perceptron as

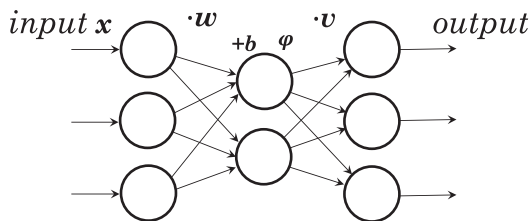


Figure 1: A schematic diagram of a three-layer perceptron presented in equation 1.1.

$$f_{(d)}(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^d v_i \varphi(\mathbf{w}_i \cdot \mathbf{x}). \quad (1.2)$$

In this article, we treat the supervised learning. The goal of the supervised learning is to find an optimal parameter  $\boldsymbol{\theta}$  so that  $f_{(d)}(\mathbf{x}; \boldsymbol{\theta})$  approximates a given target function  $T(\mathbf{x})$ . Such a problem is based on the universal approximation property of the three-layer perceptron. For a suitable activation function  $\varphi$  (e.g., sigmoidal or ReLU), the model, equation 1.2, can approximate quite a wide range of functions as the number  $d$  of hidden units tends to infinity (Cybenko, 1989; Funahashi, 1989; Sonoda & Murata, 2017).

The (averaged) gradient descent method is a standard method for finding an optimal value of  $\boldsymbol{\theta}$  numerically. Suppose that a loss function  $\ell(\mathbf{x}, \mathbf{y})$  is nonnegative and is equal to zero if and only if  $\mathbf{y} = T(\mathbf{x})$  (e.g., the squared error  $\|\mathbf{y} - T(\mathbf{x})\|^2$ ). In the gradient descent method, we aim at minimizing

$$L_{(d)}(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x}} [\ell(\mathbf{x}, f_{(d)}(\mathbf{x}; \boldsymbol{\theta}))] \quad (1.3)$$

by changing the parameter  $\boldsymbol{\theta}$  iteratively as

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \varepsilon \frac{\partial L_{(d)}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}_t), \quad (1.4)$$

where  $\varepsilon > 0$  is a learning constant. Here, we assume that the input vector  $\mathbf{x}$  is a random variable drawn according to an unknown probability distribution and  $\mathbb{E}_{\mathbf{x}}$  denotes the expectation with respect to  $\mathbf{x}$ . In order for the differential  $\partial L_{(d)} / \partial \boldsymbol{\theta}$  to make sense, we also assume that  $\ell(\mathbf{x}, \mathbf{y})$  is differentiable with respect to  $\mathbf{y}$  and that we can interchange the order of the differentiation  $\partial / \partial \boldsymbol{\theta}$  and the expectation  $\mathbb{E}_{\mathbf{x}}$ . We study the dynamical system, which represents the averaged gradient descent method with infinitesimal learning constant:

$$\frac{d\boldsymbol{\theta}}{dt} = - \frac{\partial L_{(d)}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}). \quad (1.5)$$

The parameter  $\theta$  descends along the gradient of  $L_{(d)}$  into a local minimum. In practice, the expectation in equation 1.3 is replaced with the arithmetic mean over a given data set, or, roughly, with a single realization of the random variable  $\ell(\mathbf{x}, \mathbf{f}_{(d)}(\mathbf{x}; \theta))$  for each learning iteration. Such a learning method involving some stochastic effects is called a *stochastic gradient descent method*.

Fukumizu and Amari (2000) studied singular regions arising from degeneration of hidden units of a three-layer perceptron. Here, the degeneration of hidden units means that several weight parameters  $\mathbf{w}_i$  take the same value and the effective number of hidden units becomes fewer than  $d$ . When  $m = 1$ , they found a novel type of singular region, often called a *Milnor-like attractor*. This region has both an attractive part consisting of local minima of  $L_{(d)}$  and a repulsive part consisting of saddle points. In practical learning, there may be some stochastic effects. Therefore, once the parameter  $\theta$  is trapped in the attractive part of this region, it fluctuates in the region by stochastic effects for a long time, until it reaches the repulsive part. This may cause serious stagnation of learning, called *plateau phenomena*. Later, Amari, Ozeki, Karakida, Yoshida, and Okada (2018) pointed out a notable fact that a Milnor-like attractor may not cause serious stagnation of learning when  $m \geq 2$ , which is also treated in this article.

More quantitative analyses for  $m = 1$  have also been carried out by Coussear, Ozeki, and Amari (2008), Wei, Zhang, Cousseau, Ozeki, and Amari (2008), and Amari et al. (2018) in the simplest case  $d = 2$ . In particular, Wei et al. (2008) introduced a new coordinate system in the parameter space by

$$\begin{cases} \mathbf{w} = \frac{v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2}{v_1 + v_2} \\ v = v_1 + v_2 \\ \mathbf{u} = \mathbf{w}_1 - \mathbf{w}_2 \\ z = \frac{v_1 - v_2}{v_1 + v_2} \end{cases}, \quad (1.6)$$

and claimed, based on evidence observed in numerical simulations, that when the initial point is taken near a Milnor-like attractor, the parameters  $(\mathbf{w}, v)$  quickly converge to equilibrium values  $(\mathbf{w}^*, v^*)$ . They hypothesized that this would always be the case and analyzed only the reduced dynamical system for the subparameters  $(\mathbf{u}, z)$ , setting the remaining parameters  $(\mathbf{w}, v)$  to be  $(\mathbf{w}^*, v^*)$ . However, to the best of our knowledge, no mathematical justification for this hypothesis has been established.

The objective of this article is to provide a solid ground on Amari et al.'s (2018) point of view. We introduce a new coordinate system that admits a center manifold structure around a special point on the Milnor-like attractor. By using the coordinate system, we can analyze the Milnor-like

attractor more rigorously and integrate the reduced dynamical system explicitly to obtain analytical trajectories. The obtained trajectories are comparable to the preceding work. It is confirmed by several settings of numerical simulations that trajectories in actual learning agree with the analytical ones.

In addition to the averaged gradient descent method, we also address online learning, a stochastic gradient descent method. Around a Milnor-like attractor, the behavior of sample paths by the online learning seems qualitatively different from that of trajectories by the averaged gradient descent. To investigate why they are different, we divide the dynamics of parameters into fast and slow ones, as is the case in the averaged gradient descent. In this case, we observed in numerical simulations that the fast parameters fluctuate intensively around the center manifold for the averaged system. We show that such a deviation of the fast parameter from the center manifold can influence a trend of the slow parameter.

This article is organized as follows. In section 2, we give a quick review of Amari et al.'s (2018) work. In section 3, after a brief account of the center manifold theory, we introduce a new coordinate system in the parameter space and prove that it admits the center manifold structure. In section 4, we carry out numerical simulations and observe the center manifold structure around a Milnor-like attractor. In section 5, we consider the online learning from the viewpoint of the center manifold theory. Section 6 offers concluding remarks.

## 2 Singular Region and Milnor-Like Attractor

---

In this section, we give a quick review of the Milnor-like attractor that Fukumizu and Amari (2000) found, which appears when the number  $m$  of output units is equal to 1. We also mention an interesting insight by Amari et al. (2018) for the case  $m \geq 2$ .

The parameter space of a perceptron is sometimes called a perceptron manifold. However, in many cases, it is not really a manifold since it usually contains a subset whose points correspond to the same input-output relation. Such a subset is usually referred to as a singular region. In general, there are many singular regions due to the degeneration of hidden units. For example, let us consider an  $(n-2-m)$ -perceptron. Then the subset

$$R(\mathbf{w}, \mathbf{v}) := \{\boldsymbol{\theta} = (\mathbf{w}_1, \mathbf{w}_2, \mathbf{v}_1, \mathbf{v}_2) \mid \mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}, \mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}\}$$

of the parameter space forms a typical singular region. In fact, on the subset  $R(\mathbf{w}, \mathbf{v})$ , an  $(n-2-m)$ -perceptron  $f_{(2)}(\mathbf{x}; \boldsymbol{\theta})$  is reduced to the following  $(n-1-m)$ -perceptron:

$$f_{(1)}(\mathbf{x}; \mathbf{w}, \mathbf{v}) := \mathbf{v} \varphi(\mathbf{w} \cdot \mathbf{x}).$$

On such a singular region, some properties of  $L_{(1)}$  are inherited by  $L_{(2)}$ . The following lemma implies that a criticality is a hereditary property.

**Lemma 1.** *Let  $\theta^* = (\mathbf{w}^*, \mathbf{v}^*)$  be a critical point of  $L_{(1)}$ . Then the parameter  $\theta = (\mathbf{w}_1, \mathbf{w}_2, \mathbf{v}_1, \mathbf{v}_2) = (\mathbf{w}^*, \mathbf{w}^*, \lambda \mathbf{v}^*, (1 - \lambda)\mathbf{v}^*)$  is a critical point of  $L_{(2)}$  for any  $\lambda \in \mathbb{R}$ .*

**Proof.**

$$\begin{aligned} \frac{\partial L_{(2)}}{\partial \mathbf{w}_i}(\theta) &= \mathbb{E} \left[ \left( \frac{\partial \ell(\mathbf{x}, f_{(1)}(\mathbf{x}; \theta^*))}{\partial \mathbf{y}} \cdot \mathbf{v}_i \right) \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right] = \lambda_i \frac{\partial L_{(1)}}{\partial \mathbf{w}}(\theta^*), \\ \frac{\partial L_{(2)}}{\partial \mathbf{v}_i}(\theta) &= \mathbb{E} \left[ \frac{\partial \ell(\mathbf{x}, f_{(1)}(\mathbf{x}; \theta^*))}{\partial \mathbf{y}} \varphi(\mathbf{w}^* \cdot \mathbf{x}) \right] = \frac{\partial L_{(1)}}{\partial \mathbf{v}}(\theta^*), \quad i = 1, 2, \end{aligned}$$

where  $\lambda_1 := \lambda$  and  $\lambda_2 := 1 - \lambda$ . Since  $\theta^*$  is a critical point of  $L_{(1)}$ , these are all zero. □

When  $m = 1$ , in particular, every point  $\theta \in R(\mathbf{w}^*, \mathbf{v}^*)$  is a critical point of  $L_{(2)}$ , since the parameter  $\mathbf{v}$ , as well as the output  $f^{(d)}(\mathbf{x}; \theta)$ , is a scalar quantity. In this case, the second-order property of  $L_{(1)}$  is also inherited by  $L_{(2)}$  to some extent, and the singular region  $R(\mathbf{w}^*, \mathbf{v}^*)$  may have an interesting structure, which causes serious stagnation of learning.

**Proposition 1** (Fukumizu & Amari, 2000). *Let  $m = 1$  and  $\theta^* = (\mathbf{w}^*, \mathbf{v}^*)$  be a strict local minimizer of  $L_{(1)}$  with  $\mathbf{v}^* \neq 0$ . Define an  $(n + 1) \times (n + 1)$  matrix  $H$  by*

$$H := \mathbb{E}_{\mathbf{x}} \left[ \frac{\partial \ell(\mathbf{x}, f_{(1)}(\mathbf{x}; \theta^*))}{\partial \mathbf{y}} \mathbf{v}^* \varphi''(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x} \mathbf{x}^T \right], \tag{2.1}$$

and for  $\lambda \in \mathbb{R}$ , let

$$\theta_\lambda := (\mathbf{w}^*, \mathbf{w}^*, \lambda \mathbf{v}^*, (1 - \lambda)\mathbf{v}^*).$$

*If the matrix  $H$  is positive (resp. negative) definite, then the point  $\theta = \theta_\lambda$  is a local minimizer (resp. saddle point) of  $L_{(2)}$  for any  $\lambda \in (0, 1)$  and is a saddle point (resp. local minimizer) for any  $\lambda \in \mathbb{R} \setminus [0, 1]$ . On the other hand, if the matrix  $H$  is indefinite, then the point  $\theta_\lambda$  is a saddle point of  $L_{(2)}$  for all  $\lambda \in \mathbb{R} \setminus \{0, 1\}$ .*

This proposition implies that the one-dimensional region  $R(\mathbf{w}^*, \mathbf{v}^*) = \{\theta_\lambda \mid \lambda \in \mathbb{R}\}$  may have both attractive parts and repulsive parts in the gradient descent method. Such a region is referred to as a *Milnor-like attractor* (Wei et al., 2008). The parameter  $\theta$  near the attractive part flows into the Milnor-like attractor and fluctuates in the region for a long time, until it reaches the repulsive part.

The original theorem (Fukumizu & Amari, 2000) is for an  $(n-d-1)$ -perceptron that contains  $(n-(d-1)-1)$ -perceptron as a subnetwork and that

the phenomenon itself is universal with respect to the number  $d$  of hidden units. The proposition above for  $(n-2-1)$ -perceptron is a minimal version.

We also remark that the point  $\theta_\lambda$  cannot be a strict local minimizer since  $L_{(2)}$  takes the same value on the singular region  $\{\theta_\lambda \mid \lambda \in \mathbb{R}\}$  and is flat along its direction. The proof of proposition 1 is given mainly by a discussion of the Hessian matrix of  $L_{(2)}$ ; however, we need to treat higher-order derivatives of  $L_{(2)}$  than the second order, since the Hessian matrix degenerates on the singular region (see appendix A).

Let us suppose a situation where a three-layer perceptron has some redundant hidden units to represent the target function  $T(x)$ . Mathematically, we suppose that a true parameter  $\theta_{true}$  exists ( $T(x) = f_{(2)}(x; \theta_{true})$ ) and that it lies in the singular region  $R(\mathbf{w}^*, v^*)$ . In this case, the function  $L_{(2)}$  takes the same value  $L_{(1)}(\mathbf{w}^*, v^*) = 0$  on  $R(\mathbf{w}^*, v^*)$ . Therefore, every point of  $R(\mathbf{w}^*, v^*)$  becomes a global minimizer of  $L_{(2)}$ , and a Milnor-like attractor does not appear. In fact, one can check that the assumption of proposition 1 fails as follows. For each  $x \in \mathbb{R}^n$ , we obtain

$$\partial \ell(x, f_{(1)}(x; \mathbf{w}^*, v^*)) / \partial \mathbf{y} = 0,$$

since a loss function  $\ell(x, \mathbf{y})$  takes its minimum 0 at  $\mathbf{y} = T(x) = f_{(1)}(x; \mathbf{w}^*, v^*)$ . This implies that the matrix  $H$  becomes the zero matrix. Thus,  $H$  is neither positive nor negative definite.

We next treat the case when  $m \geq 2$ . There also exists a one-dimensional region consisting of critical points due to lemma 1. However, in this case, the region becomes simply repulsive and does not have an attractive part, as the following theorem asserts.

**Theorem 1.** *Let  $\theta^* = (\mathbf{w}^*, v^*)$  be a local minimizer of  $L_{(1)}$ . If the  $m \times (n + 1)$  matrix*

$$\mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(1)}(x; \theta^*))}{\partial \mathbf{y}} \varphi'(\mathbf{w}^* \cdot x) x^T \right] \tag{2.2}$$

*is nonzero, then  $\theta_\lambda = (\mathbf{w}^*, \mathbf{w}^*, \lambda v^*, (1 - \lambda)v^*)$  is a saddle point of  $L_{(2)}$  for any  $\lambda \in \mathbb{R}$ , where we regard the derivative  $\partial \ell / \partial \mathbf{y}$  as a column vector.*

Amari et al. (2018) stated a prototype of theorem 1, although they did not give a full proof. In fact, we found that some additional assumption was necessary to prove their assertion. In theorem 1, we have added a mild assumption that the matrix, equation 2.2, is nonzero. Note that since  $\theta^* = (\mathbf{w}^*, v^*)$  is a local minimizer of  $L_{(1)}$ , it holds that

$$\mathbf{0} = \frac{\partial L_{(1)}}{\partial \mathbf{w}}(\boldsymbol{\theta}^*) = (\mathbf{v}^*)^T \mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(1)}(x; \boldsymbol{\theta}^*))}{\partial \mathbf{y}} \boldsymbol{\varphi}'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right].$$

Thus, the matrix 2.2 has a kernel whose dimension is greater than or equal to one. Hence, the assumption automatically fails when  $m = 1$ . This is an underlying mechanism for proposition 1.

### 3 Center Manifold of Milnor-Like Attractor

In their analysis of an  $(n-2-1)$ -perceptron, Wei et al. (2008) introduced a coordinate transformation,

$$\left\{ \begin{array}{l} \mathbf{w} = \frac{v_1 \mathbf{w}_1 + v_2 \mathbf{w}_2}{v_1 + v_2} \\ v = v_1 + v_2 \\ \mathbf{u} = \mathbf{w}_1 - \mathbf{w}_2 \\ z = \frac{v_1 - v_2}{v_1 + v_2} \end{array} \right. , \tag{3.1}$$

and claimed that the parameters  $(\mathbf{w}, v)$  quickly converge to  $(\mathbf{w}^*, v^*)$  when the initial point is taken near a Milnor-like attractor. Amari et al. (2018) mentioned that the dynamics in this coordinate system should be analyzed by using the center manifold theory, and they analyzed only the reduced dynamical system for the subparameters  $(\mathbf{u}, z)$ , setting the remaining parameters  $(\mathbf{w}, v)$  to be  $(\mathbf{w}^*, v^*)$ . Strictly speaking, however, their coordinate system does not admit any center manifold structure, and their claim is at the stage of hypothesis.

In this section, we give a rigorous justification for their hypothesis. We first give a quick review of the center manifold theory and then introduce a new coordinate system under which the center manifold structures arise near certain points on the Milnor-like attractor.

**3.1 Brief Review of Center Manifold.** Suppose that we are given a dynamical system,

$$\left\{ \begin{array}{l} \dot{\mathbf{x}}(t) = A\mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), \mathbf{y}(t)) \\ \dot{\mathbf{y}}(t) = B\mathbf{y}(t) + \mathbf{g}(\mathbf{x}(t), \mathbf{y}(t)) \end{array} \right. , \tag{3.2}$$

for the parameters  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$ , where  $A$  and  $B$  are constant matrices and  $\mathbf{f}$  and  $\mathbf{g}$  are  $C^2$  functions such that they, along with their first derivatives, vanish at the origin. We assume that all the eigenvalues of  $A$  have zero real parts, while all the eigenvalues of  $B$  have negative real parts. This assumption means that parameter  $\mathbf{y}$  converges to the origin exponentially quickly, and the parameter  $\mathbf{x}$  is driven only by the higher-order terms of  $\mathbf{f}$

and evolves very slowly compared with  $y$ . Since  $f$  and  $g$  are of the second order with respect to  $x$  and  $y$ , the assumptions for equation 3.2 imply that the coefficient matrix of the linearization of the system has the form

$$\begin{pmatrix} A & O \\ O & B \end{pmatrix}.$$

**Definition 1.** A set  $S \subset \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$  is said to be a local invariant manifold of equation 3.2 if for  $(x_0, y_0) \in S$ , the solution  $(x(t), y(t))$  of equation 3.2 with  $(x(0), y(0)) = (x_0, y_0)$  is in  $S$  for  $|t| < T$  with some  $T > 0$ .

**Definition 2.** A local invariant manifold represented in the form of  $y = h(x)$  is called a local center manifold (or simply a center manifold) if  $h$  is differentiable and satisfies  $h(0) = 0$  and  $\frac{\partial h}{\partial x}(0) = O$ .

The following center manifold theorems give us a method of simplifying a dynamical system around an equilibrium point.

**Proposition 2** (Center manifold theorem 1: Carr, 1981). Equation 3.2 has a center manifold  $y = h(x)$  for  $\|x\| < \delta$ , for some  $\delta > 0$  and  $C^2$  function  $h$ .

**Proposition 3** (Center manifold theorem 2: Carr, 1981). Suppose that the origin  $u = 0$  is a stable equilibrium point of the reduced dynamical system

$$\dot{u}(t) = Au(t) + f(u(t), h(u(t))). \tag{3.3}$$

Let  $(x(t), y(t))$  be a solution of equation 3.2 with the initial value  $(x_0, y_0)$ . Then if  $\|(x_0, y_0)\|$  is sufficiently small, there exists a solution  $u(t)$  of equation 3.3 such that

$$\begin{aligned} x(t) &= u(t) + O(e^{-\gamma t}), \\ y(t) &= h(u(t)) + O(e^{-\gamma t}), \end{aligned}$$

as  $t \rightarrow \infty$ , where  $\gamma$  is a positive constant.

Proposition 3 asserts that the parameter  $(x, y)$  approaches the center manifold  $y = h(x)$  quickly and then evolves along it. Thus, the dynamical system, equation 3.2, around the origin is essentially controlled by the slow parameter  $x$  and reduced to the lower-dimensional system.

**3.2 Main Results.** Let us return to the analysis of an  $(n-2-1)$ -perceptron. In a column vector representation, the dynamical system, equation 1.5, for the  $(n-2-1)$ -perceptron is written as

$$\dot{\theta} = - \left( \frac{\partial L_{(2)}(\theta)}{\partial \theta} \right)^T.$$



We consider another coordinate system  $\xi = \xi(\theta)$  and investigate the dynamical system in it. By the coordinate transformation, the dynamical system above is transformed to

$$\dot{\xi} = -\frac{\partial \xi}{\partial \theta} \left( \frac{\partial \xi}{\partial \theta} \right)^T \left( \frac{\partial L_{(2)}}{\partial \xi}(\xi) \right)^T. \tag{3.4}$$

Thus, the coefficient matrix of its linearization at a critical point  $\xi = \xi^*$  is

$$\begin{aligned} \frac{\partial \dot{\xi}}{\partial \xi}(\xi^*) &= -\frac{\partial}{\partial \xi} \left\{ \frac{\partial \xi}{\partial \theta} \left( \frac{\partial \xi}{\partial \theta} \right)^T \left( \frac{\partial L_{(2)}}{\partial \xi}(\xi) \right)^T \right\} \Bigg|_{\xi=\xi^*} \\ &= -\frac{\partial \xi}{\partial \theta} \left( \frac{\partial \xi}{\partial \theta} \right)^T \frac{\partial^2 L_{(2)}}{\partial \xi \partial \xi}(\xi^*), \end{aligned}$$

where we used  $(\partial L_{(2)}/\partial \xi)(\xi^*) = \mathbf{0}$ . This relation implies that the coefficient matrix has the same rank as the Hessian matrix  $(\partial^2 L_{(2)}/\partial \xi \partial \xi)(\xi^*)$ . In particular, the rank of the coefficient matrix of the linearization does not depend on the choice of a coordinate system.

We focus on the dynamics of the learning process around the two points,  $\theta = \theta_0, \theta_1$ , which are boundaries of the repulsive and attractive parts of a Milnor-like attractor  $\{\theta_\lambda | \lambda \in \mathbb{R}\}$ . Concretely, they are denoted as  $\theta_0 = (\mathbf{w}^*, \mathbf{w}^*, 0, v^*)$  and  $\theta_1 = (\mathbf{w}^*, \mathbf{w}^*, v^*, 0)$ , where  $(\mathbf{w}^*, v^*)$  is a minimizer of the loss  $L_{(1)}$  for the  $(n-1)$ -perceptron as mentioned in proposition 1. This is because the rank of the Hessian matrix at  $\theta_\lambda$  degenerates by one dimension for  $\lambda \neq 0, 1$  and by  $n + 2$  dimension for  $\lambda = 0, 1$ , which is shown in appendix A.

We introduce a new coordinate system  $\xi = (\mathbf{w}, v, \mathbf{u}, z)$  by

$$\begin{cases} \mathbf{w} = \frac{v_1(\mathbf{w}_1 - \mathbf{w}^*) + v_2(\mathbf{w}_2 - \mathbf{w}^*)}{v^*} + \mathbf{w}^* \\ v = v_1 + v_2 \\ \mathbf{u} = \frac{v_2(\mathbf{w}_1 - \mathbf{w}^*) - v_1(\mathbf{w}_2 - \mathbf{w}^*)}{v^*} \\ z = v_1 - v_2 \end{cases} \tag{3.5}$$

This formula defines a coordinate system on the region  $\{v_1^2 + v_2^2 \neq 0\}$ . Under the coordinate system, equation 3.5, the critical points  $\theta_\lambda$  are denoted as

$$\xi_\lambda = (\mathbf{w}^*, v^*, \mathbf{0}, (2\lambda - 1)v^*).$$

In particular,  $\xi_0 = (\mathbf{w}^*, v^*, \mathbf{0}, -v^*)$  and  $\xi_1 = (\mathbf{w}^*, v^*, \mathbf{0}, v^*)$ . Now we arrive at the main theorem of this article:

**Theorem 2.** *In the coordinate system  $\xi = (\mathbf{w}, v, \mathbf{u}, z)$ , the dynamical system, equation 1.5, admits a center manifold structure around the critical points  $\theta = \theta_0, \theta_1$  in which  $(\mathbf{w}, v)$  converge exponentially fast.*

To prove the theorem, we make use of the following lemma.

**Lemma 2.** *If the matrix  $X$  is positive definite and  $Y$  is positive semidefinite, all the eigenvalues of the matrix  $XY$  are nonnegative.*

**Proof.** The matrix  $XY$  is rewritten as

$$XY = X^{\frac{1}{2}}(X^{\frac{1}{2}}Y X^{\frac{1}{2}})X^{-\frac{1}{2}},$$

where  $X^{\frac{1}{2}}$  is a unique positive-definite matrix such that  $(X^{\frac{1}{2}})^2 = X$ . Here, the matrix  $Z := X^{\frac{1}{2}}Y X^{\frac{1}{2}}$  is positive semidefinite. Hence, for each eigenvector  $\mathbf{a}$  of  $Z$ , the vector  $X^{\frac{1}{2}}\mathbf{a}$  is an eigenvector of the matrix  $XY$ , and the corresponding eigenvalue is nonnegative.  $\square$

**Proof of Theorem 2.** The proof is essentially based on a straightforward calculation. The coefficient matrix of the linearization of the dynamical system, equation 3.4, under the coordinate system, equation 3.5, splits into  $(\mathbf{w}, v)$  part and  $(\mathbf{u}, z)$  part for  $\lambda = 0, 1$ . In fact, for  $\lambda \in \mathbb{R}$ , the negative of the coefficient matrix is written as

$$\tilde{A}_\lambda = \begin{matrix} & \overbrace{\hspace{10em}}^{\mathbf{w}, v} & & \overbrace{\hspace{5em}}^{\mathbf{u}} & \overbrace{\hspace{2em}}^z \\ \mathbf{w}, v \left\{ \right. & \left( (1 + 2k_\lambda)Q + \frac{1+3k_\lambda}{1+2k_\lambda}H \right. & (1 + 2k_\lambda)P & -\frac{(-1+2\lambda)k_\lambda}{1+2k_\lambda}H & 0 \\ & 2P^T & 2R & 0 & 0 \\ \mathbf{u} \left\{ \right. & -\frac{(-1+2\lambda)k_\lambda}{1+2k_\lambda}H & 0 & -\frac{k_\lambda}{1+2k_\lambda}H & 0 \\ z \left\{ \right. & 0 & 0 & 0 & 0 \end{matrix} \right\},$$

and the system is written as

$$\dot{\xi} = -\tilde{A}_\lambda(\xi - \xi_\lambda) + \tilde{g}_\lambda(\xi),$$

where  $\tilde{g}_\lambda$  is the higher-order term, which vanish at the  $\xi = \xi_\lambda$  together with its first derivative. Here,

$$\begin{aligned} k_\lambda &:= (1 - \lambda)\lambda, \\ P &:= \mathbb{E}_x [(\partial^2 \ell) v^* \varphi(\mathbf{w}^* \cdot \mathbf{x}) \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}], \\ Q &:= \mathbb{E}_x [(\partial^2 \ell) (v^*)^2 \varphi'(\mathbf{w}^* \cdot \mathbf{x})^2 \mathbf{x} \mathbf{x}^T], \\ R &:= \mathbb{E}_x [(\partial^2 \ell) \varphi(\mathbf{w}^* \cdot \mathbf{x})^2], \end{aligned}$$

$$\begin{aligned} \partial \ell &:= \frac{\partial \ell(x, f_{(1)}(x; \theta^*))}{\partial y}, \\ \partial^2 \ell &:= \frac{\partial^2 \ell(x, f_{(1)}(x; \theta^*))}{\partial y^2}, \end{aligned}$$

and  $H$  is the matrix defined by equation 2.1.  $H$  and  $Q$  are matrices,  $P$  is a column vector, and  $R$  is a scalar. For  $\lambda = 0$  or  $1$ , the negative of the coefficient matrix is of the form

$$\tilde{A}_0 = \tilde{A}_1 = \begin{matrix} & \overbrace{\begin{matrix} w, v \end{matrix}} & \overbrace{\begin{matrix} u \end{matrix}} & \overbrace{\begin{matrix} z \end{matrix}} \\ \begin{matrix} w, v \\ u \\ z \end{matrix} \left\{ \begin{matrix} Q + H & P & 0 & 0 \\ 2P^T & 2R & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \right. \end{matrix} \quad (3.6)$$

We show that all the eigenvalues of  $(w, v)$ -block of the coefficient matrix  $-\tilde{A}_0$  are strictly negative. Recall that the coefficient matrix at a critical point  $\xi^*$  of the dynamical system, equation 1.5, is given by

$$-\frac{\partial \xi}{\partial \theta} \left( \frac{\partial \xi}{\partial \theta} \right)^T \frac{\partial^2 L_{(2)}}{\partial \xi \partial \xi}(\xi^*).$$

Applying lemma 2 to  $X = (\partial \xi / \partial \theta)(\partial \xi / \partial \theta)^T$  and  $Y = (\partial^2 L_{(2)} / \partial \xi \partial \xi)(\xi_0)$ , all the eigenvalues of the coefficient matrix  $-\tilde{A}_0$  are nonpositive. One can see that the Hessian matrix is positive semidefinite and degenerates by  $n + 2$  dimension at  $\theta = \theta_0, \theta_1$  in appendix A. Since a coordinate transformation preserves the rank of the coefficient matrix of linearization,  $\tilde{A}_0$  degenerates by  $n + 2$  dimension, which is equal to the size of  $(u, z)$ -block. This implies that the  $(w, v)$ -block is of full rank, and thus all the eigenvalues of  $(w, v)$ -block are strictly negative. It is proved similarly for  $\xi = \xi_1$ .

Due to proposition 2, there are center manifolds parametrized by  $(u, z)$  around  $\theta = \theta_0, \theta_1$  respectively. □

**3.3 Reduced Dynamical System.** By virtue of proposition 3 and theorem 2, we can assume that the dynamics of the gradient descent is on the center manifold near the points  $\theta = \theta_0, \theta_1$ . Thus, we can reduce the dynamical system into that of  $(u, z)$ . Recalling the coefficient matrix, equation 3.6, we can see that  $\dot{u}$  and  $\dot{z}$  have no first-order terms. In more detail, calculating the Taylor expansion of  $(\dot{u}, \dot{z})$  up to the second order around  $\xi = \xi_1$ , we obtain

$$\begin{aligned} \dot{u} = & \frac{1}{v^*} \left\{ -(P \cdot (w - w^*))(u + (w - w^*)) \right. \\ & - (v - v^*)(RI + \frac{1}{2}H)(u + (w - w^*)) \\ & \left. + \frac{1}{2}(z - v^*)H(u + (w - w^*)) \right\} + O(\|\xi - \xi_1\|^3), \end{aligned} \tag{3.7}$$

$$\begin{aligned} \dot{z} = & \frac{1}{v^*} \left\{ -(w - w^*)^T Q(u + (w - w^*)) \right. \\ & - (v - v^*)(P \cdot (u + (w - w^*))) \\ & \left. - \frac{1}{2}(w - w^*)^T H(w - w^*) + \frac{1}{2}u^T H u \right\} + O(\|\xi - \xi_1\|^3), \end{aligned} \tag{3.8}$$

where  $I$  denotes the  $(n + 1) \times (n + 1)$  identity matrix. Now we consider the reduced dynamical system on the center manifold. Here, the center manifold  $(w, v) = h(u, z)$  satisfies that

$$h(u, z) = \begin{bmatrix} w(u, z) \\ v(u, z) \end{bmatrix} = \begin{bmatrix} w^* \\ v^* \end{bmatrix} + O(\|u, z - v^*\|^2),$$

by definition. This gives an approximation of the dynamics on the center manifold near  $\xi = \xi_1$  as

$$\begin{aligned} \dot{u} &= \frac{1}{2v^*} (z - v^*)Hu + O(\|u, z - v^*\|^3), \\ \dot{z} &= \frac{1}{2v^*} u^T H u + O(\|u, z - v^*\|^3). \end{aligned} \tag{3.9}$$

Neglecting the higher-order terms, we can integrate this equation to obtain

$$\|u\|^2 = (z - v^*)^2 + C, \tag{3.10}$$

where  $C$  is an integral constant.

Around the point  $\xi = \xi_0$ , we obtain the similar dynamics,

$$\begin{aligned} \dot{u} &= -\frac{1}{2v^*} (z + v^*)Hu + O(\|u, z + v^*\|^3), \\ \dot{z} &= -\frac{1}{2v^*} u^T H u + O(\|u, z + v^*\|^3), \end{aligned}$$

and the relation,

$$\|u\|^2 = (z + v^*)^2 + C.$$

We remark that theorem 2 is valid even when there exists a true parameter in the singular region  $R(\mathbf{w}^*, v^*)$ ; however, in this case, such a simple form of the reduced dynamical system as equation 3.9 is not obtained. As mentioned above, this case implies that  $H$  becomes the zero matrix. Then, the second-order terms of the reduced dynamical system, equation 3.9, vanish, and the third-order terms become dominant. Thus, we have to take into account the cross terms between  $(\mathbf{w} - \mathbf{w}^*, v - v^*)$  and  $(\mathbf{u}, z - v^*)$ . It needs to calculate the center manifold  $(\mathbf{w}, v) = \mathbf{h}(\mathbf{u}, z)$  up to the second order, which makes the analysis complicated.

Finally, we remark on a difference between our analysis and previous work. Wei et al. (2008) have studied a reduced dynamical system in the vicinity of the whole part of a Milnor-like attractor. On the other hand, a center manifold is defined locally, and center manifolds around each of two points cannot be connected at a midpoint in general. Thus, one cannot discuss a center manifold defined around the entire region of a Milnor-like attractor.

**3.4 More General Models.** Our results can be extended to a more general model including multilayer perceptrons whose output is one-dimensional. In this section, we consider a parameterized family of functions that can be written as

$$\begin{aligned}
 f(\mathbf{x}; \boldsymbol{\theta}) &:= g(\mathbf{x}, \boldsymbol{\tau}) + v_1 \phi(\mathbf{x}; \mathbf{w}_1, \boldsymbol{\tau}) + v_2 \phi(\mathbf{x}; \mathbf{w}_2, \boldsymbol{\tau}), \\
 \boldsymbol{\theta} &= (\mathbf{w}_1, \mathbf{w}_2, v_1, v_2, \boldsymbol{\tau}),
 \end{aligned}
 \tag{3.11}$$

where we assume that  $g(\mathbf{x}; \boldsymbol{\tau})$  and  $\phi(\mathbf{x}; \mathbf{w}, \boldsymbol{\tau})$  are twice differentiable with respect to  $\boldsymbol{\tau}$  and  $(\mathbf{w}, \boldsymbol{\tau})$ , respectively.

A multilayer perceptron  $\tilde{\alpha}^{(L)}(\mathbf{x}; \boldsymbol{\theta})$  with  $L$  layers defined recursively as

$$\begin{aligned}
 \alpha^{(0)}(\mathbf{x}; \boldsymbol{\theta}) &:= (1, x_1, \dots, x_{n_0})^T, \\
 \tilde{\alpha}^{(\ell)}(\mathbf{x}; \boldsymbol{\theta}) &:= W^{(\ell)} \alpha^{(\ell-1)}(\mathbf{x}; \boldsymbol{\theta}), \\
 \alpha^{(\ell)}(\mathbf{x}; \boldsymbol{\theta}) &:= \left( 1, \varphi(\tilde{\alpha}_1^{(\ell)}(\mathbf{x}; \boldsymbol{\theta})), \dots, \varphi(\tilde{\alpha}_{n_\ell}^{(\ell)}(\mathbf{x}; \boldsymbol{\theta})) \right)^T, \quad 1 \leq \ell \leq L, \\
 \boldsymbol{\theta} &= (W^{(1)}, \dots, W^{(L)}),
 \end{aligned}$$

for each  $\mathbf{x} = (x_1, \dots, x_{n_0})^T \in \mathbb{R}^{n_0}$ , where  $W^{(\ell)}$  is an  $n_\ell \times (n_{\ell-1} + 1)$  matrix,  $n_0, n_1, \dots, n_L \in \mathbb{N}$ , and  $\varphi$  is a twice-differentiable activation function. Assuming that the dimension  $n_L$  of the output is equal to 1 and denoting

$$\begin{aligned}
 W^{(L-1)} &= [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \cdots \quad \mathbf{w}_{n_{L-1}}]^T, \\
 W^{(L)} &= [v_0 \quad v_1 \quad v_2 \quad \cdots \quad v_{n_{L-1}}],
 \end{aligned}$$

a multilayer perceptron is represented as model 3.11 by letting

$$\begin{aligned} \phi(\mathbf{x}; \mathbf{w}, \boldsymbol{\tau}) &:= \varphi(\mathbf{w} \cdot \alpha^{(L-2)}(\mathbf{x}; \boldsymbol{\tau})), \\ g(\mathbf{x}; \boldsymbol{\tau}) &:= \sum_{i=3}^{n_{L-1}} v_i \varphi(\mathbf{w}_i \cdot \alpha^{(L-2)}(\mathbf{x}; \boldsymbol{\tau})), \\ \boldsymbol{\tau} &:= (W^{(1)}, \dots, W^{(L-2)}, \mathbf{w}_3, \dots, \mathbf{w}_{n_{L-1}}, v_0, v_3, \dots, v_{n_{L-1}}). \end{aligned}$$

Our main result is extended to the model, equation 3.11 as follows. Let  $\boldsymbol{\theta}^* = (\mathbf{w}^*, v^*, \boldsymbol{\tau}^*)$  be a strict local minimizer of an averaged loss function for the degenerate model:

$$f_{(1)}(\mathbf{x}; \tilde{\mathbf{w}}, \tilde{v}, \boldsymbol{\tau}) = g(\mathbf{x}; \boldsymbol{\tau}) + \tilde{v} \phi(\mathbf{x}; \tilde{\mathbf{w}}, \boldsymbol{\tau}).$$

Then the coordinate system  $\boldsymbol{\xi} = (\mathbf{w}, v, \boldsymbol{\tau}, \mathbf{u}, z)$ , given by formula 3.5, admits a center manifold structure around the two points  $\boldsymbol{\theta} = (\mathbf{w}^*, \mathbf{w}^*, 0, v^*, \boldsymbol{\tau}^*)$ ,  $(\mathbf{w}^*, \mathbf{w}^*, v^*, 0, \boldsymbol{\tau}^*)$ , and the dynamical system is reduced to that of  $(\mathbf{u}, z)$ . This is confirmed by the argument of the coefficient matrix of the linearization, similar to theorem 2.

#### 4 Numerical Simulations

---

In the previous section, we showed that the dynamics of  $(\mathbf{w}, v)$  are fast and those of  $(\mathbf{u}, z)$  are slow under the coordinate system, equation 3.5. In this section, we verify this fact by numerical simulations.

**4.1 Example 1.** As the first example, we set the input dimension to be  $n = 1$  and choose the teacher function  $T : \mathbb{R} \rightarrow \mathbb{R}$  defined by

$$T(x) := 2 \tanh(x) - \tanh(4x),$$

where  $\tanh$  is the hyperbolic tangent function. The shape of  $T$  is shown in Figure 2 by the solid black line. We set the activation function  $\varphi$  as  $\tanh$ . Thus, the target function  $T$  can be represented by the (1-2-1)-perceptron with no bias terms:

$$f_{(2)}(x; \boldsymbol{\theta}) = v_1 \varphi(w_1 x) + v_2 \varphi(w_2 x),$$

and the true parameter is  $(w_1, w_2, v_1, v_2) = (1, 4, 2, -1)$ . We also discard the bias terms  $w_1^0$  and  $w_2^0$  of the student perceptron. This makes the matrix  $H$  scalar valued, and hence the assumption of proposition 1 holds trivially.

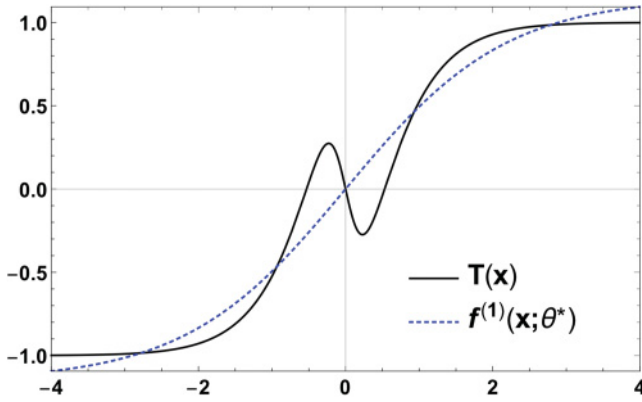


Figure 2: The target function  $T(x)$  and the (1-1-1)-perceptron  $f_{(1)}(x; \theta^*)$ , which corresponds to the local minimizer  $\theta^*$ .

We assume that the data set  $\{x_s\}_{s=1}^S$  is given and that the probability distribution of  $x$  is the empirical distribution on that data set. Then the transition formula, equation 1.4, of the parameter  $\theta$  is rewritten as

$$\theta^{(t+1)} = \theta^{(t)} - \varepsilon \frac{1}{S} \sum_{s=1}^S \left. \frac{\partial \ell(x_s, f_2(x_s; \theta))}{\partial \theta} \right|_{\theta = \theta^{(t)}}.$$

In this simulation, we set the size  $S$  of the data set to be 1000, and data  $\{x_s\}_{s=1}^S$  are drawn identical and independently distributed (i.i.d.) according to  $N(0, 2^2)$ . Here,  $N(\mu, \sigma^2)$  denotes the gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .

For a data set given as above, we obtained a local minimizer  $\theta^* = (w^*, v^*) = (0.459, 1.15)$  of  $L_{(1)}$ . The shape of the function that corresponds to the local minimizer is shown by the dashed blue line in Figure 2. The value of  $H$  is approximately 0.0472. Since  $H > 0$ , the attractive region is  $\{\theta_\lambda \mid \lambda \in (0, 1)\}$ , due to proposition 1.

Figure 3 displays time evolutions of each parameter in the first 1500 iterations from 50 different initial points. We chose an initial parameter  $\theta^{(0)} = (w_1^{(0)}, w_2^{(0)}, v_1^{(0)}, v_2^{(0)})$  by

$$\begin{aligned} w_1^{(0)} &= w^* + \zeta_1, & w_2^{(0)} &= w^* + \zeta_2, \\ v_1^{(0)} &= v^* + \frac{1}{2}(\zeta_3 + \zeta_4), & v_2^{(0)} &= \frac{1}{2}(\zeta_3 - \zeta_4), \end{aligned}$$

so that  $v = v^* + \zeta_3$  and  $z = v^* + \zeta_4$ , where  $\zeta_1, \zeta_2 \sim U(-0.2, 0.2)$ , and  $\zeta_3, \zeta_4 \sim U(-0.2, 0.2)$ . Here,  $U(a, b)$  denotes the uniform distribution on the interval

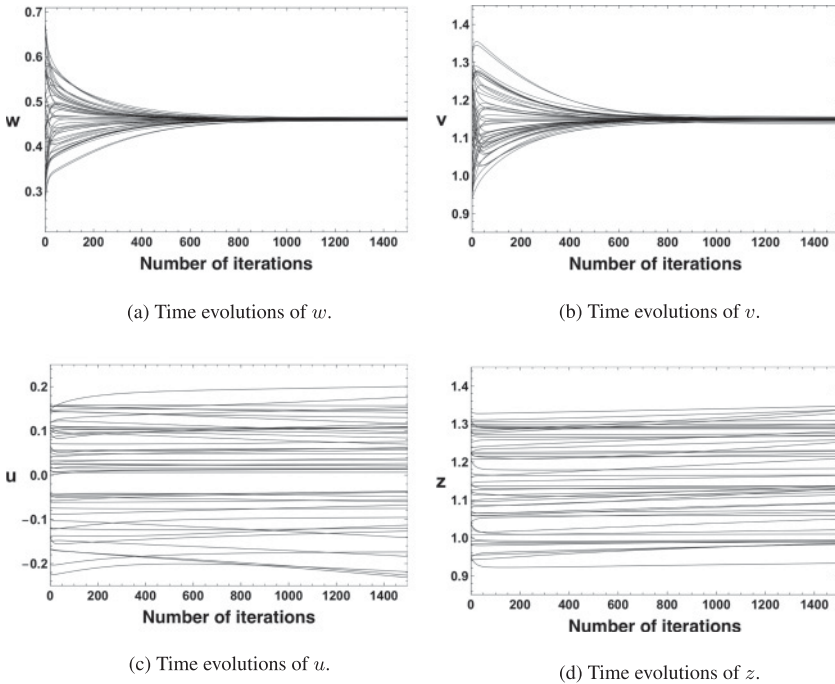


Figure 3: Time evolutions of each parameter for the first 1500 iterations. Each trajectory of  $w$  or  $v$  quickly converges to the equilibrium point  $w^* = 0.459$ ,  $v^* = 1.15$ , respectively. Trajectories of  $u$  and  $z$  evolve very slowly compared with  $w$  and  $v$ .

$[a, b] \subset \mathbb{R}$ . We set the learning rate  $\varepsilon$  to be 0.05 and the number of iterations to be 20,000. We can see that the parameters  $w$  and  $v$  converge to their equilibriums exponentially fast (see Figures 3a and 3b), while  $u$  and  $z$  evolve slowly (see Figures 3c and 3d).

Figure 4 shows evolutions on the  $(z, \|u\|^2)$ -plane. The red circles in the figure represent initial points. When  $w = w^*$  and  $v = v^*$ , the  $z$ -axis is a Milnor-like attractor, and the region  $|z| < v^*$  is the attractive part of it. We can check that parameters near the attractive region are trapped and those near the repulsive region are escaping. The intersection point of the line  $z = v^*$  and  $z$ -axis corresponds to the point  $\theta = \theta_1$ , the boundary of the attractive and repulsive parts of the Milnor-like attractor. The analytical trajectories, equation 3.10, are plotted as dashed blue curves. Numerical evolutions of the parameter follow the analytical trajectories considerably well around  $\theta = \theta_1$ . We see in the figure that some instances of time evolutions change their direction sharply. This is because the fast dynamics of  $w$  and  $v$  are the main dynamics at the beginning of the learning, while the



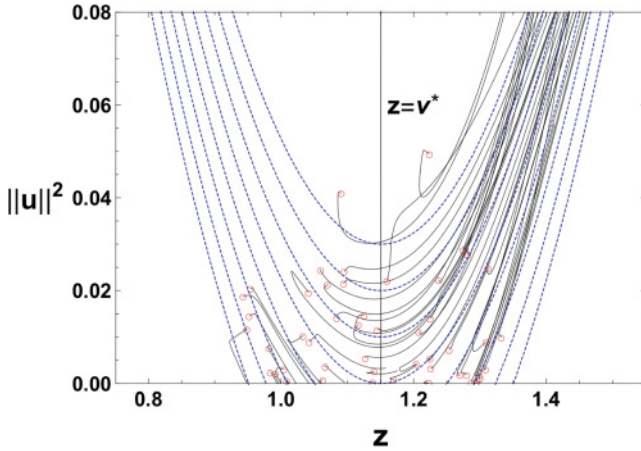


Figure 4: Trajectories on the  $(z, \|u\|^2)$ -plane obtained by learning for 20,000 iterations (solid black curves) and analytical trajectories (dashed blue curves) near  $\theta = \theta_1 = (w^*, w^*, v^*, 0)$ . Red circles represent initial points.

slow dynamics of  $u$  and  $z$  become dominant after  $w$  and  $v$  converge to the center manifold.

**4.2 Example 2.** As the second example, we consider a three-layer perceptron whose input dimension is  $n = 2$ . Let the teacher function  $T : \mathbb{R}^2 \rightarrow \mathbb{R}$  be given by

$$T(x_1, x_2) := 0.75 \operatorname{Sgm}(2.5x_1 - 2.5x_2) + \operatorname{Sgm}(2.5x_1 + 2.5x_2) + 0.5,$$

where  $\operatorname{Sgm}$  is the logistic sigmoidal function, which is defined by

$$\operatorname{Sgm}(x) := \frac{1}{1 + e^{-2x}} = \frac{1}{2}(1 + \tanh(x)).$$

Figure 5 shows the shape of the teacher function  $T(x)$ . We use a perceptron with no bias terms also in this simulation and choose  $\operatorname{Sgm}$  as the activation function. Note that a (2-2-1)-perceptron is unable to represent the target function in this case.

Also in this simulation, we assume that the probability distribution of  $x$  is the empirical distribution on a fixed data set  $\{x_s\}_{s=1}^S$ . We set the number  $S$  of the data set to be 1000, and draw a data set  $\{x_s\}_{s=1}^S$  i.i.d. according to  $N(0, I_2)$ , where  $I_2$  denotes the  $2 \times 2$  identity matrix. We chose a realization  $\{x_s\}_{s=1}^S$  as above and obtained a local minimizer  $\theta^* = (w^*, v^*)$  of  $L_{(1)}$ , where  $w^* = (0.399, 0.0652)$  and  $v^* = 2.76$ . Figure 6 shows the shape of the

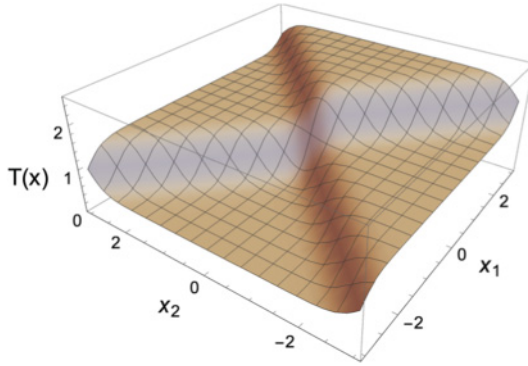


Figure 5: The target teacher function  $T(x)$ .

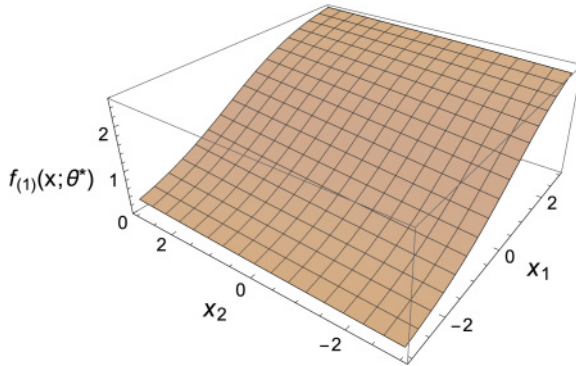


Figure 6: The (2-1-1)-perceptron  $f_{(1)}(x; \theta^*)$  that corresponds to the local minimizer  $\theta^*$ .

(2-1-1)-perceptron corresponding to the local minimizer. The matrix  $H$  is numerically computed as

$$\begin{pmatrix} -0.044 & -0.026 \\ -0.026 & -0.20 \end{pmatrix}.$$

Since this  $H$  is negative definite, the attractive region is  $\{\theta_\lambda \mid \lambda \in \mathbb{R} \setminus [0, 1]\}$  due to proposition 1.

Figures 7a to 7d show time evolutions of each parameter in the first 500 iterations from 50 different initial points. We chose initial parameters of the (2-2-1)-perceptron as

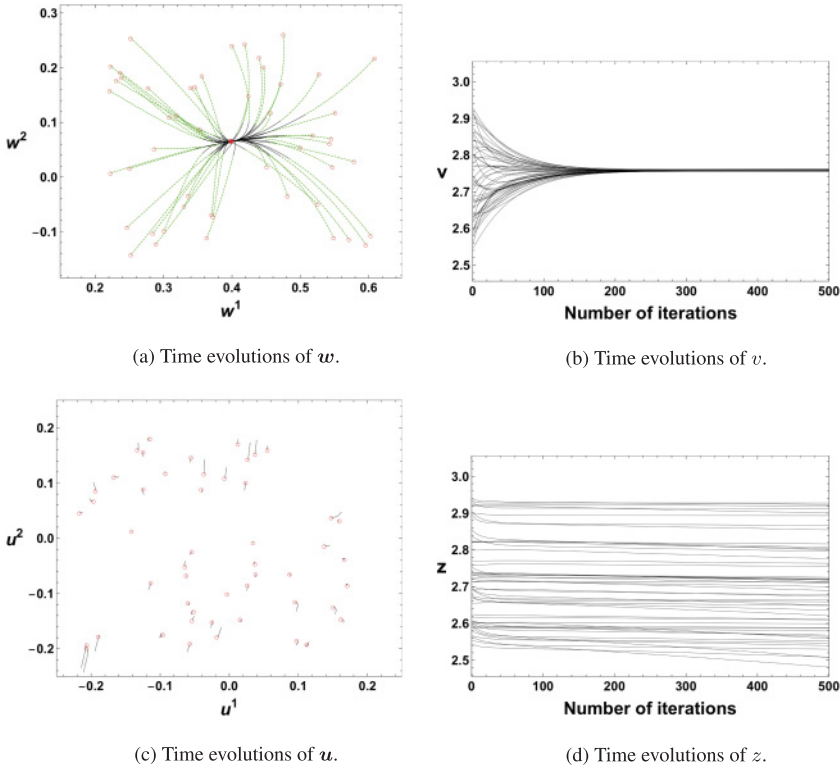


Figure 7: Time evolutions of each parameter for the first 500 iterations. Each trajectory of  $w$  or  $v$  quickly converges to the equilibrium point  $w^* = (0.399, 0.0652)$ ,  $v^* = 2.76$  respectively. However, trajectories of  $u$  and  $z$  evolve very slowly compared with  $w$  and  $v$ . In panel a, the first 30 iterations are dashed in green to display the speed of convergence, and the red point at the center represents  $w = w^*$ . Red circles in panels a and c represent initial points.

$$w_1^{(0)} = w^* + \zeta_1, \quad w_2^{(0)} = w^* + \zeta_2,$$

$$v_1^{(0)} = v^* + \frac{1}{2}(\zeta_3 + \zeta_4), \quad v_2^{(0)} = \frac{1}{2}(\zeta_3 - \zeta_4),$$

where  $\zeta_1, \zeta_2 \sim U(-0.2, 0.2)^2$  and  $\zeta_3, \zeta_4 \sim U(-0.2, 0.2)$ . We set the learning rate  $\varepsilon$  to be 0.05 and the number of iterations to be 20,000. In this simulation, since  $w$  and  $u$  are two-dimensional, their evolutions are not displayed as time series but as trajectories on each plane. The red circles in Figures 7a and 7c represent initial values of  $w$  and  $u$ , respectively. Figures 7a and 7b show that the parameters  $w$  and  $v$  converge to their equilibrium  $w^*$  and  $v^*$  very quickly. To display how quick the convergence is, the first 30 iterations

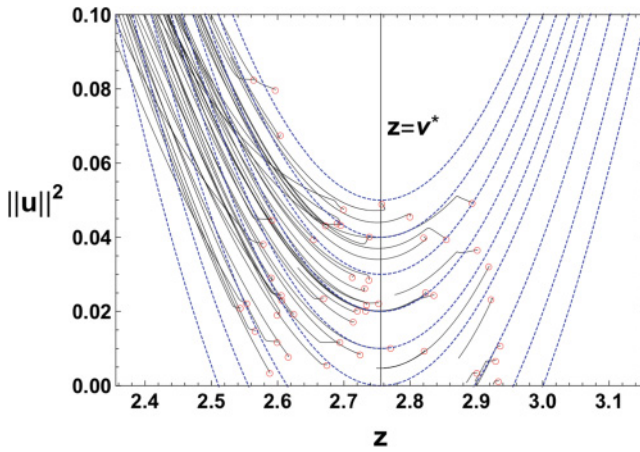


Figure 8: Trajectories on the  $(z, \|\mathbf{u}\|^2)$ -plane obtained by learning for 20,000 iterations (solid black curves) and analytical trajectories near  $\theta = \theta_1$  (dashed blue curves). Red circles represent the initial points.

are dashed in green in Figure 7a. In contrast, Figures 7c and 7d show that the parameters  $\mathbf{u}$  and  $z$  evolve very slowly.

Figure 8 plots time evolutions of the parameter  $\theta$  on the  $(z, \|\mathbf{u}\|^2)$ -plane, which means the plane whose axes indicate the values of  $z$  and  $\|\mathbf{u}\|^2$ . We can check that parameters near the attractive part  $\{|z| > v^*, \|\mathbf{u}\|^2 = 0\}$  of the Milnor-like attractor are trapped and that those near the repulsive part  $\{|z| < v^*, \|\mathbf{u}\|^2 = 0\}$  are escaping. The numerical evolutions follow the analytical flows (dashed blue curves) well also in this case.

## 5 Aspects of Online Learning

In this section, we discuss stochastic effects in the learning process. Thus far, we have analyzed the dynamical system, equation 1.5, driven by the averaged gradient. In practice, the averaged gradient is estimated by the arithmetic mean of the instantaneous gradient  $\partial \ell(\mathbf{x}, f_{(d)}(\mathbf{x}; \theta)) / \partial \theta$  over a large number of input data. However, taking the arithmetic mean for each update of the parameter demands high computational cost. In order to reduce the cost, the expectation is often replaced by a single realization of the instantaneous gradient. Such a method is called *online learning*, a typical stochastic gradient descent method. Mathematically, it is given by

$$\theta^{(t+1)} = \theta^{(t)} - \varepsilon \left. \frac{\partial}{\partial \theta} \ell(\mathbf{x}_t, f_{(d)}(\mathbf{x}_t; \theta)) \right|_{\theta=\theta^{(t)}}, \quad (5.1)$$

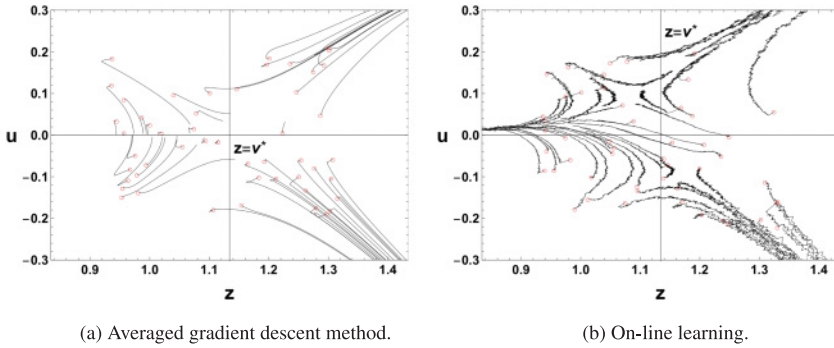


Figure 9: Trajectories on the  $(z, u)$ -plane obtained by the averaged gradient descent method, equation 1.5, and the online learning, equation 5.1, for 20,000 iterations. Red circles represent initial values.

where  $\{x_t\}_t$  are i.i.d. realizations of the input data  $x$ . Unlike the deterministic dynamical system, equation 1.5, the system 5.1 is a randomized dynamical system.

In numerical simulations, we found that sample paths of the online learning seem quite different from trajectories obtained in the averaged gradient descent method. We set the distribution over which the loss function  $L_{(1)}$  is averaged to be  $N(0, 2^2)$  and obtained a local minimizer of  $L_{(1)}$  as  $\theta^* = (w^*, v^*) \approx (0.472, 1.13)$ . We carried out numerical simulations of the online learning in the same setting as example 1 in section 4. Figure 9a shows numerical trajectories of the averaged gradient descent on the  $(z, u)$ -plane around  $\theta = \theta_1$ . In order to approximate the averaged gradient descent sufficiently, we used the empirical distribution on a data set of 10,000 data drawn i.i.d. according to  $N(0, 2^2)$ . Figure 9b shows sample paths of the online learning for a common input data sequence  $\{x_t\}_t$ . In contrast to the averaged gradient descent, in the online learning, some sample paths move from region  $\{|z| > v^*\}$  to  $\{|z| < v^*\}$ . Such sample paths are observed even when we use another realization of the input data sequence, and its dynamics seems qualitatively different from the averaged one.

In order to investigate this phenomenon, we observe the evolution of the parameters, again in the coordinate system, equation 3.5. Figures 10a to 10d show time evolutions of each parameter in the first 1500 iterations of the online learning. The parameters  $(w, v)$  evolve very fast compared with  $(u, z)$  also in this case. However, in this case,  $(w, v)$  does not converge to its equilibrium point  $(w^*, v^*) \approx (0.472, 1.13)$ , but fluctuates stochastically around  $(w^*, v^*)$ .

Based on these observations, we suppose that  $w$  and  $v$  run over a sufficiently wide range of their values to be integrated, while  $u$  and  $z$  move in a small range. Then we assume that the dynamics of  $(u, z)$  is integrated

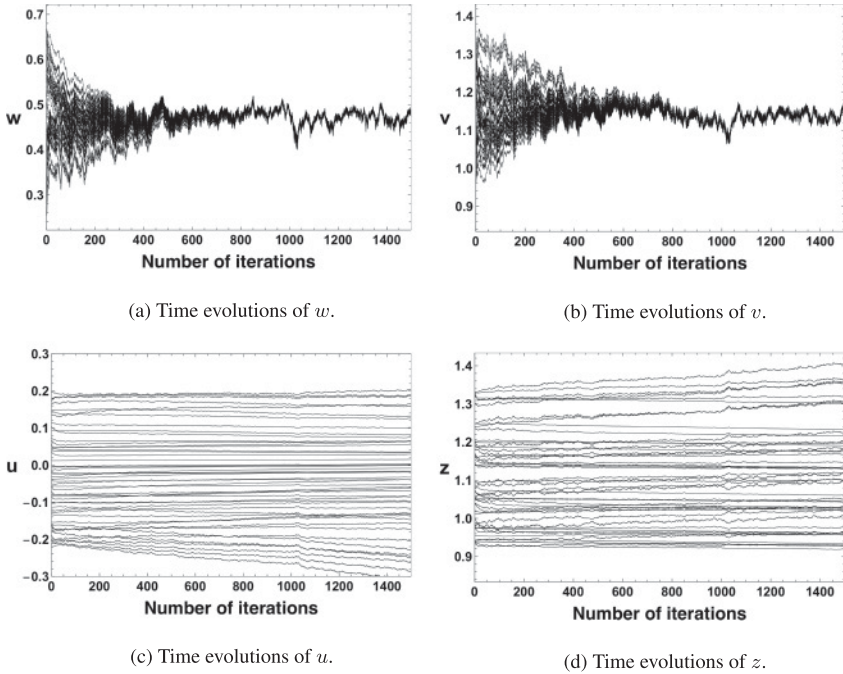


Figure 10: Time evolutions of each parameter for the first 1500 iterations of the online learning. Each trajectory of  $w$  or  $v$  fluctuates intensively around the equilibrium point  $w^* \approx 0.472$ ,  $v^* \approx 1.134$ , respectively. Trajectories of  $u$  and  $z$  evolve very slowly compared with  $w$  and  $v$  in this case.

with respect to  $(w, v)$  according to some stationary distribution. We further assume that  $(w, v)$  are distributed around  $(w^*, v^*)$  with finite variance. By integrating the Taylor expansions, equations 3.7 and 3.8, with  $(w, v)$ , we obtain the following dynamical system near  $\theta = \theta_1$ :

$$\begin{aligned} \dot{u} &= \frac{1}{2} (z - v^*)Hu + C_1, \\ \dot{z} &= \frac{1}{2} u^T Hu + C_2. \end{aligned} \quad (5.2)$$

Here,  $C_1$  and  $C_2$  are constants resulting from the variance and covariance of  $(w, v)$ . Figure 11 shows the analytical trajectories of the dynamical system, equation 5.2, where  $C_1 = 1.71 \times 10^{-4}$  and  $C_2 = -3.06 \times 10^{-4}$  are determined heuristically. One can find that the deterministic dynamical system, equation 5.2, gives similar trajectories to sample paths of the online learning presented in Figure 9b.

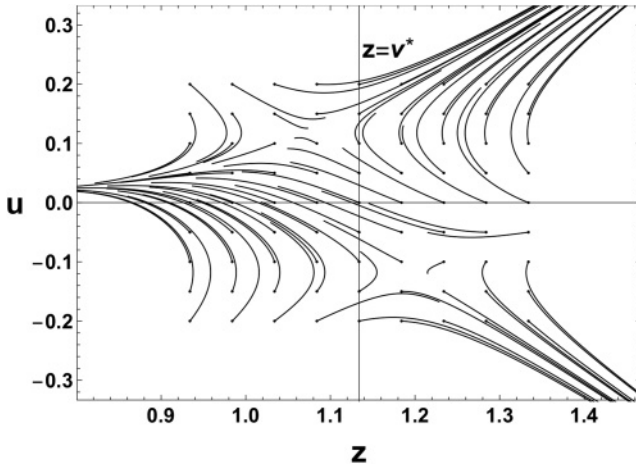


Figure 11: Analytical trajectories on the  $(z, u)$ -plane given by the dynamical system 5.2 with  $C_1 = 1.71 \times 10^{-4}$  and  $C_2 = -3.06 \times 10^{-4}$ .

We deduce that a fluctuation of the parameter around a center manifold causes constants  $C_1$  and  $C_2$  working as drift terms and that it makes the dynamics of the online learning qualitatively different from those of the averaged gradient descent. This example suggests that stochastic effects can influence a macroscopic flow of the learning process via a center manifold structure.

## 6 Conclusion

In this article, we first gave a quick review of a mechanism that causes plateau phenomena in a three-layer perceptron—in particular, how degeneration of hidden units gives rise to a Milnor-like attractor consisting of both attractive and repulsive parts. We next investigated the dynamics of learning around special points on a Milnor-like attractor and proved the existence of the center manifold. We also succeeded in integrating the reduced dynamical system to obtain an analytical form of a trajectory. We performed several numerical simulations to demonstrate the accuracy of our results. As an application of the center manifold structure, we gave an explanation for a characteristic behavior of the online learning.

Unfortunately, the two examples presented in section 4 were the only ones that we could find in which the assumptions of proposition 1 are fulfilled. This might suggest that the appearance of a Milnor-like attractor would be a rather rare situation in a perceptron that has bias terms. In fact, just by replacing the activation function  $\text{Sgm}$  with  $\tanh$  in example 2, the matrix  $H$  becomes indefinite and the assumption of proposition 1

is violated. Finding more suggestive examples that shed light on the complex behavior of the dynamics of learning is an important subject for future study.

In section 5, we investigated stochastic effects of the online learning from an intermediate viewpoint between fully stochastic and averaged dynamics. We made use of the center manifold of the averaged dynamics and discussed an integration with quickly fluctuating parameters. There have been many reports of qualitative differences between stochastic and deterministic methods; however, there are few general theories for analyzing such dissimilarities. We expect that the intermediate viewpoint in this article can be a clue to clarify stochastic effects in learning.

**Appendix: Proofs of Proposition 1 and Theorem 1** \_\_\_\_\_

This appendix gives proofs of proposition 1 and Theorem 1. The proof of proposition 1 is based on the analysis of the Hessian matrix of  $L_{(2)}$ . However, the Hessian at the point  $\theta_\lambda$  becomes singular, since  $L_{(2)}(\theta_\lambda)$  is constant along  $\lambda \in \mathbb{R}$ . Thus, we need to take into account higher-order derivatives of  $L_{(2)}$ , which is overlooked in the Fukumizu and Amari (2000). The prototype of theorem 1 was given by Amari et al. (2018); however, they proved it only for a special case. Here, we give a rigorous proof with an additional mild assumption.

**Proof of Proposition 1.** We introduce a new coordinate system  $\xi = (w, v, u, z)$  by

$$\begin{cases} w = \frac{v_1 w_1 + v_2 w_2}{v_1 + v_2} \\ v = v_1 + v_2 \\ u = w_1 - w_2 \\ z = v_1 - v_2 \end{cases}, \tag{A.1}$$

where  $v_1 + v_2 \neq 0$ . Under this coordinate system, the point  $\theta_\lambda$  is denoted as  $\xi_\lambda = (w^*, v^*, \mathbf{0}, (2\lambda - 1)v^*)$ . Note that each point  $\xi = \xi_\lambda$  is a critical point of  $L_{(2)}(\xi)$  by lemma 1. The inverse transformation is given as

$$\begin{cases} w_1 = w + \frac{v - z}{2v} u \\ w_2 = w - \frac{v + z}{2v} u \\ v_1 = \frac{v + z}{2} \\ v_2 = \frac{v - z}{2} \end{cases}$$



We observe the Hessian matrix  $\text{Hess}(\xi_\lambda)$  of  $L_{(2)}(\xi)$  at  $\xi_\lambda$  for each  $\lambda \in \mathbb{R}$ . For all  $\xi$  such that  $u = 0$  (or, equivalently,  $w_1 = w_2$ ), using the inverse transformation formula above,

$$\begin{aligned} \frac{\partial L_{(2)}}{\partial u}(\xi) &= \mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(2)}(x; \xi))}{\partial y} \frac{v^2 - z^2}{4v} (\varphi'(w_1 \cdot x) - \varphi'(w_2 \cdot x))x \right] \\ &= 0, \\ \frac{\partial L_{(2)}}{\partial z}(\xi) &= -\mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(2)}(x; \xi))}{\partial y} \left( \frac{v+z}{4v} \varphi'(w_1 \cdot x) - \frac{v-z}{4v} \varphi'(w_2 \cdot x) \right) (u \cdot x) \right] \\ &\quad + \mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(2)}(x; \xi))}{\partial y} \frac{1}{2} (\varphi(w_1 \cdot x) - \varphi(w_2 \cdot x)) \right] \\ &= 0. \end{aligned}$$

Here, we left  $w_1$  and  $w_2$  for notational simplicity. We then have

$$\begin{aligned} \frac{\partial^2 L_{(2)}}{\partial \gamma \partial z}(\xi_\lambda) &= 0, \\ \frac{\partial^2 L_{(2)}}{\partial \gamma \partial u}(\xi_\lambda) &= 0, \end{aligned}$$

where  $\gamma = w, v, z$ . Hence, the matrix  $\text{Hess}(\xi_\lambda)$  has the form

$$\begin{matrix} & \begin{matrix} \overbrace{w, v} & \overbrace{u} & \overbrace{z} \end{matrix} \\ \begin{matrix} w, v \\ u \\ z \end{matrix} \left\{ \begin{matrix} * & * & 0 & 0 \\ * & * & 0 & 0 \\ 0 & 0 & * & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \right. \end{matrix}.$$

The  $(w, v)$ -block is equal to the Hessian matrix of  $L_{(1)}$  at  $\theta^*$  and is positive definite. In fact, for any  $\xi$  such that  $u = 0$ , noting that  $f_{(2)}(x; \xi) \equiv f_{(1)}(x; w, v)$ ,

$$\begin{aligned} \frac{\partial L_{(2)}}{\partial w}(\xi) &= \mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(1)}(x; w, v))}{\partial y} v \varphi'(w \cdot x)x \right] = \frac{\partial L_{(1)}}{\partial w_{(1)}}(w, v), \\ \frac{\partial L_{(2)}}{\partial v}(\xi) &= \mathbb{E}_x \left[ \frac{\partial \ell(x, f_{(1)}(x; w, v))}{\partial y} \varphi(w \cdot x) \right] = \frac{\partial L_{(1)}}{\partial v_{(1)}}(w, v). \end{aligned}$$

The  $(w, v)$ -block of  $\text{Hess}(\xi_\lambda)$  is given by differentiating the equations above with  $(w, v)$  and thus equal to the Hessian matrix  $(\partial^2 L_{(1)} / \partial \theta_{(1)} \partial \theta_{(1)})(\theta^*)$ , which is positive definite since  $\theta_{(1)} = \theta^*$  is a strict local minimizer of  $L_{(1)}(\theta_{(1)})$ .

On the other hand, we have

$$\begin{aligned} \frac{\partial^2 L_{(2)}(\xi_\lambda)}{\partial \mathbf{u} \partial \mathbf{u}} &= \mathbb{E}_x \left[ \frac{\partial \ell(\mathbf{x}, f_{(2)}(\mathbf{x}; \xi_\lambda))}{\partial y} \lambda(1 - \lambda) v^* \varphi''(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x} \mathbf{x}^T \right] \\ &= \lambda(1 - \lambda)H, \end{aligned}$$

where  $H$  is the matrix defined in equation 2.1. Thus, when  $H$  is positive definite, the Hessian matrix  $\text{Hess}(\xi_\lambda)$  is positive semidefinite if and only if  $\lambda \in [0, 1]$ . This proves that  $\xi_\lambda$  is a saddle point of  $L_{(2)}(\xi)$  for  $\lambda \in \mathbb{R} \setminus [0, 1]$ . Similarly, when  $H$  is negative definite,  $\xi_\lambda$  is a saddle point of  $L_{(2)}(\xi)$  for  $\lambda \in (0, 1)$ . When  $H$  is indefinite, so is  $\text{Hess}(\xi_\lambda)$  for every  $\lambda \in \mathbb{R} \setminus \{0, 1\}$ .

We show that the point  $\xi = \xi_\lambda$  is a local minimizer of  $L_{(2)}(\xi)$  for  $\lambda \in (0, 1)$  when the matrix  $H$  is positive definite. For  $\lambda \in (0, 1)$ , since the matrix  $\text{Hess}(\xi_\lambda)$  is positive semidefinite, the Taylor series of  $(L_{(2)}(\xi) - L_{(2)}(\xi_\lambda))$  up to the second order is nonnegative in a neighborhood of  $\xi = \xi_\lambda$ . However, higher-order terms may influence the sign of the Taylor series, since the coefficients of the terms of  $z^2$  and  $\gamma z$  are zero for  $\gamma = \mathbf{w}, v, \mathbf{u}$ . Since the terms of  $\gamma^k z^\ell$  with  $k \geq 2$  are dominated by the term  $\gamma^2$  near  $\xi = \xi_\lambda$ , we check if the coefficients of the terms of  $z^\ell$  and  $\gamma z^\ell$  are equal to zero for all  $\ell \geq 1$ . For all  $\xi$  such that  $\mathbf{u} = \mathbf{0}$ ,  $f_{(2)}(\mathbf{x}; \xi) = v\varphi(\mathbf{w} \cdot \mathbf{x})$  is a constant as a function of  $z$ , and so is  $L_{(2)}(\xi)$ . Hence, we obtain

$$\begin{aligned} \frac{\partial^\ell L_{(2)}(\xi_\lambda)}{\partial z^\ell} &= 0, \\ \frac{\partial^{\ell+1} L_{(2)}(\xi_\lambda)}{\partial \gamma \partial z^\ell} &= 0, \quad \ell \geq 1, \end{aligned}$$

for  $\gamma = \mathbf{w}, v$ . Therefore, it suffices to check if

$$\frac{\partial^{\ell+1} L_{(2)}(\xi_\lambda)}{\partial \mathbf{u} \partial z^\ell} = 0, \quad \ell \geq 1.$$

Since we have already seen that  $(\partial L_{(2)} / \partial \mathbf{u})(\xi) = 0$  for all  $\xi$  such that  $\mathbf{u} = \mathbf{0}$ , this equality is confirmed. Finally, we have shown that the coefficients of higher-order terms are all zero, and that they do not influence the sign of the Taylor series. This shows that  $(L_{(2)}(\xi) - L_{(2)}(\xi_\lambda))$  is nonnegative near  $\xi = \xi_\lambda$ , and hence the proof is complete. In the case that the matrix  $H$  is negative, it is proved similarly.  $\square$

**Proof of Theorem 1.** Since  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are no longer scalars, we cannot use the coordinate system, given by equation A.1. Therefore, in order to analyze

the Hessian matrix, we introduce another coordinate system  $\xi = (\mathbf{w}, \mathbf{v}, \mathbf{u}, \mathbf{z})$  as

$$\begin{cases} \mathbf{w} = \frac{\mathbf{w}_1 + \mathbf{w}_2}{2} \\ \mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2 \\ \mathbf{u} = \frac{\mathbf{w}_1 - \mathbf{w}_2}{2} \\ \mathbf{z} = \mathbf{v}_1 - \mathbf{v}_2 \end{cases}.$$

In this coordinate system, the point  $\theta_\lambda$  is denoted as  $\xi_\lambda = (\mathbf{w}^*, \mathbf{v}^*, \mathbf{0}, (2\lambda - 1)\mathbf{v}^*)$ .

Fix  $\lambda \in \mathbb{R}$  arbitrarily. We show that the Hessian matrix  $\text{Hess}(\xi_\lambda)$  of  $L_{(2)}(\xi)$  at  $\xi = \xi_\lambda$  has both positive and negative eigenvalues. It suffices to show that the  $(\mathbf{w}, \mathbf{v})$ -part of the Hessian is positive definite and that the  $(\mathbf{u}, \mathbf{z})$ -part is not positive semidefinite. They imply that the full Hessian matrix  $\text{Hess}(\xi_\lambda)$  is neither negative semidefinite nor positive semidefinite, and hence that  $\text{Hess}(\xi_\lambda)$  is indefinite.

One can check that the  $(\mathbf{w}, \mathbf{v})$ -part of  $\text{Hess}(\xi_\lambda)$  is equal to the Hessian matrix of  $L_{(1)}$  at  $\theta^*$  by direct calculation. Since  $\theta^*$  is a strict local minimizer of  $L_{(1)}$ , this is positive definite.

On the other hand, we have

$$\frac{\partial^2 L_{(2)}}{\partial \mathbf{z} \partial \mathbf{u}}(\xi_\lambda) = \mathbb{E}_x \left[ \frac{\partial \ell(\mathbf{x}, f_{(2)}(\mathbf{x}; \xi_\lambda))}{\partial \mathbf{y}} \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right], \tag{A.2}$$

$$\frac{\partial^2 L_{(2)}}{\partial \mathbf{z} \partial \mathbf{z}}(\xi_\lambda) = O, \tag{A.3}$$

and thus the  $(\mathbf{u}, \mathbf{z})$ -part  $B$  of  $\text{Hess}(\xi_\lambda)$  is written as

$$B = \begin{pmatrix} \frac{\partial^2 L_{(2)}}{\partial \mathbf{u} \partial \mathbf{u}}(\xi_\lambda) & \mathbb{E}_x \left[ \frac{\partial \ell(\mathbf{x}, f_{(2)}(\mathbf{x}; \xi_\lambda))}{\partial \mathbf{y}} \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right]^T \\ \mathbb{E}_x \left[ \frac{\partial \ell(\mathbf{x}, f_{(2)}(\mathbf{x}; \xi_\lambda))}{\partial \mathbf{y}} \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right] & O \end{pmatrix},$$

where we treat  $\partial \ell / \partial \mathbf{y}$  as a column vector. Here, since  $f_{(2)}(\mathbf{x}; \xi_\lambda) \equiv f_{(1)}(\mathbf{x}; \theta^*)$ , the nondiagonal block, equation A.2, is equal to the matrix, equation 2.2, and thus is nonzero by assumption. Hence, the  $(\mathbf{u}, \mathbf{z})$ -part above is not positive semidefinite. In fact, choosing vectors  $\mathbf{a} \in \mathbb{R}^m$  and  $\mathbf{b} \in \mathbb{R}^{n+1}$  such that

$$\mathbf{a}^T \mathbb{E}_x \left[ \frac{\partial \ell(\mathbf{x}, f_{(2)}(\mathbf{x}; \xi_\lambda))}{\partial \mathbf{y}} \varphi'(\mathbf{w}^* \cdot \mathbf{x}) \mathbf{x}^T \right] \mathbf{b} < 0,$$

the vector  $\mathbf{c}_\varepsilon := (\varepsilon \mathbf{b}^T, \mathbf{a}^T)^T$  satisfies  $\mathbf{c}_\varepsilon^T B \mathbf{c}_\varepsilon < 0$  for sufficiently small  $\varepsilon > 0$ . Such vectors  $(\mathbf{a}, \mathbf{b})$  always exist—for instance, by letting  $a_i = -\rho$ ,  $b_j = 1$ ,

and other entries be 0, where  $(i, j)$  is an index such that  $(i, j)$ -entry  $\rho$  of the matrix  $A.2$  is nonzero.  $\square$

## Acknowledgments

---

I express my gratitude to Akio Fujiwara for his helpful guidance, discussions, and advice. I thank as well Yuzuru Sato for many discussions and helpful comments. And I thank the anonymous referees for their insightful suggestions to improve this article.

## References

---

- Amari, S.-I., Ozeki, T., Karakida, R., Yoshida, Y., & Okada, M. (2018). Dynamics of learning in MLP: Natural gradient and singularity revisited. *Neural Computation*, 30(1), 1–33.
- Carr, J. (1981). *Applications of centre manifold theory*. New York: Springer.
- Cousseau, F., Ozeki, T., & Amari, S.-I. (2008). Dynamics of learning in multilayer perceptrons near singularities. *IEEE Transactions on Neural Networks*, 19(8), 1313–1328.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2(4), 303–314.
- Fukumizu, K., & Amari, S.-I. (2000). Local minima and plateaus in hierarchical structures of multilayer perceptrons. *Neural Networks*, 13(3), 317–327.
- Funahashi, K. (1989). On the approximate realization of continuous mappings by neural networks. *Neural Networks*, 2(3), 183–192.
- Sonoda, S., & Murata, N. (2017). Neural network with unbounded activation functions is universal approximator. *Applied and Computational Harmonic Analysis*, 43(2), 233–268.
- Wei, H., Zhang, J., Cousseau, F., Ozeki, T., & Amari, S.-I. (2008). Dynamics of learning near singularities in layered networks. *Neural Computation*, 20(3), 813–843.

---

Received July 26, 2019; accepted November 26, 2019.