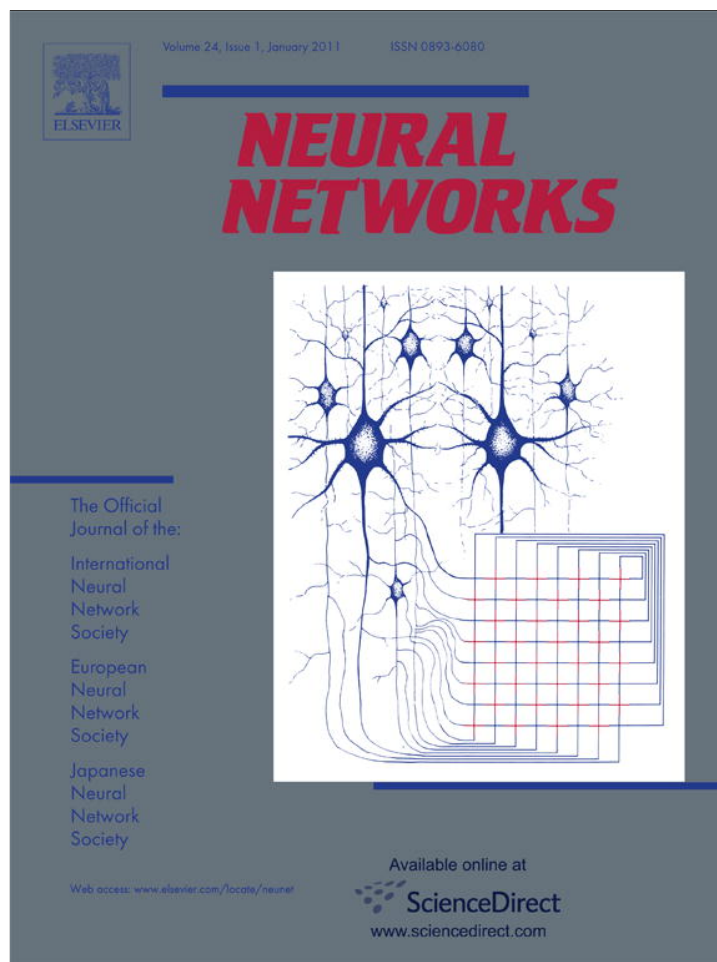


Provided for non-commercial research and education use.  
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

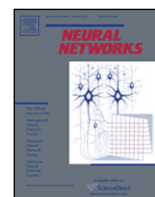
In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at ScienceDirect

## Neural Networks

journal homepage: [www.elsevier.com/locate/neunet](http://www.elsevier.com/locate/neunet)Convergence analysis of online gradient method for BP neural networks<sup>☆</sup>Wei Wu<sup>a,\*</sup>, Jian Wang<sup>a,b</sup>, Mingsong Cheng<sup>a</sup>, Zhengxue Li<sup>a</sup><sup>a</sup> School of Mathematical Sciences, Dalian University of Technology, Dalian, 116024, PR China<sup>b</sup> School of Mathematics and Computational Sciences, Petroleum University of China, Dongying, 257061, PR China

## ARTICLE INFO

## Article history:

Received 28 March 2010

Received in revised form 8 September 2010

Accepted 9 September 2010

## Keywords:

Neural networks

Backpropagation learning

Online gradient method

Weak convergence

Strong convergence

## ABSTRACT

This paper considers a class of online gradient learning methods for backpropagation (BP) neural networks with a single hidden layer. We assume that in each training cycle, each sample in the training set is supplied in a stochastic order to the network exactly once. It is interesting that these stochastic learning methods can be shown to be deterministically convergent. This paper presents some weak and strong convergence results for the learning methods, indicating that the gradient of the error function goes to zero and the weight sequence goes to a fixed point, respectively. The conditions on the activation function and the learning rate to guarantee the convergence are relaxed compared with the existing results. Our convergence results are valid for not only S–S type neural networks (both the output and hidden neurons are Sigmoid functions), but also for P–P, P–S and S–P type neural networks, where S and P represent Sigmoid and polynomial functions, respectively.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

Artificial neural network has been a hot topic in recent years in cognitive science, computational intelligence and intelligent information processing. Backpropagation (BP) is the most broadly used learning method for feedforward neural networks. It was first proposed by Werbos (1974) in his Ph.D. thesis, and has been rediscovered several times (LeCun, 1985; Parker, 1982; Rumelhart, Hinton, & Williams, 1986). There are two practical ways to implement the backpropagation algorithm: batch updating approach and online updating approach. Corresponding to the standard gradient method, the batch updating approach accumulates the weight correction over all the training samples before actually performing the update. On the other hand, the online updating approach updates the network weights immediately after each training sample is fed. Some authors compare the two different training schemes for feedforward neural networks (Heskes & Wiegierinck, 1996; Nakama, 2009; Wilson & Martinez, 2003). Heskes and Wiegierinck (1996) reveal several asymptotic properties of the two schemes. Wilson and Martinez (2003) explain why batch training is almost always slower than online training (often orders of magnitude slower) especially on large training sets. Nakama (2009) theoretically analyzes the convergence properties of the two schemes applied to quadratic

loss functions and shows the exact degrees to which the training set size, the variance of the per-instance gradient, and the learning rate affect the rate of convergence for each scheme.

There are three approaches for online training of BP neural networks according to different fashions of sampling. The first approach is OGM-CS (completely stochastic order): At each learning step, one of the samples is drawn at random from the training set and presented to the network (Finnoff, 1994; Heskes & Wiegierinck, 1996; Terence, 1989; Wilson & Martinez, 2003). The second approach is OGM-SS (special stochastic order): In each training cycle, each sample in the training set is supplied in a stochastic order to the network exactly once (Heskes & Wiegierinck, 1996; Li & Ding, 2005; Li, Wu, & Tian, 2004; Nakama, 2009). The third approach is OGM-F (fixed order): In each training cycle, each sample in the training set is supplied in a fixed order to the network exactly once (Heskes & Wiegierinck, 1996; Mangasarian & Solodov, 1994; Wu & Xu, 2002; Wu, Feng, Li, & Xu, 2005; Xu, Zhang, & Jin, 2009).

Naturally, the existing convergence results for OGM-CS are mostly asymptotic convergence with a probabilistic nature as the size of training samples goes to infinity (Bertsekas & Tsitsiklis, 1996; Chakraborty & Pal, 2003; Fine & Mukherjee, 1999; Finnoff, 1994; Liang, Feng, Lee, Lim, & Lee, 2002; Tadic & Stankovic, 2000; Terence, 1989; Zhang, Wu, Liu, & Yao, 2009). Deterministic convergence can be obtained for OGM-SS and OGM-F (Li et al., 2004; Mangasarian & Solodov, 1994; Shao, Wu, & Liu, 2007; Wu & Xu, 2002; Wu et al., 2005; Wu, Feng, & Li, 2002; Wu & Shao, 2003; Wu, Shao, & Qu, 2005; Xu et al., 2009). It is interesting to see that the learning method OGM-SS with stochastic nature enjoys deterministic convergence. The convergence result is a

<sup>☆</sup> Project supported by the National Natural Science Foundation of China (No. 10871220).

\* Corresponding author.

E-mail address: [wuweiw@dlut.edu.cn](mailto:wuweiw@dlut.edu.cn) (W. Wu).

bit easier to prove for OGM-F than for OGM-SS. But we have reason to believe, and our experience shows, that OGM-SS behaves numerically better than OGM-F since the stochastic nature of the learning procedure survives in OGM-SS (Li & Ding, 2005; Li et al., 2004).

To guarantee the convergence, it is commonly required that the learning rate  $\eta_m$  satisfies the assumptions  $\sum_{m=1}^{\infty} \eta_m = \infty$  and  $\sum_{m=1}^{\infty} \eta_m^2 < \infty$  as in Bertsekas and Tsitsiklis (1996) and Tadic and Stankovic (2000) for OGM-CS. An extra assumption  $\lim_{m \rightarrow \infty} \eta_m / \eta_{m+1} = 1$  was introduced by Xu et al. (2009) for OGM-F. A special condition which is basically  $\eta_m = O(1/m)$  was required in Li et al. (2004), Shao et al. (2007), Wu and Xu (2002), Wu et al. (2005), Wu et al. (2002), Wu and Shao (2003) and Wu et al. (2005) for OGM-F and OGM-SS.

To obtain the strong convergence result, which means that the weight sequence converges to a fixed point, Wu et al. (2005) introduced an additional assumption: the number of the stationary points of the error function is finite. A more relaxed condition is used in Xu et al. (2009): the gradient of the error function has at most countably infinite number of stationary points.

The aim of this paper is to present a comprehensive study on the weak and strong convergence for OGM-F and OGM-SS, indicating that the gradient of the error function goes to zero and the weight sequence goes to a fixed point, respectively. These convergence results improve the existing results in Li et al. (2004), Shao et al. (2007), Wu and Xu (2002), Wu et al. (2005), Wu et al. (2002), Wu and Shao (2003), Wu et al. (2005) and Xu et al. (2009) such that the conditions on the activation function and the learning rate to guarantee the convergence are much relaxed. Specifically, we make the following contributions:

- The extra condition  $\lim_{m \rightarrow \infty} \eta_m / \eta_{m+1} = 1$  for the learning rate is removed which is a requisite in Xu et al. (2009).
- The convergence results are valid for both OGM-F and OGM-SS.
- The convergence results apply not only to S–S type neural networks (both the output and hidden neurons are Sigmoid functions), but also to P–P, P–S and S–P type neural networks, where S and P represent Sigmoid and polynomial functions, respectively.
- The restrictive assumptions for the strong convergence in Wu et al. (2005) and Xu et al. (2009) are relaxed such that the stationary points set of the error function is only required not to contain any interior point.
- We assume that the derivative  $g'$  of the activation function is Lipschitz continuous on any bounded closed interval. This improves the corresponding conditions in Wu et al. (2005), which require the boundedness of the second derivative  $g''$ , and in Xu et al. (2009), which require  $g'$  to be Lipschitz continuous and uniformly bounded on the whole  $R$ .

Let us make a few remarks on the above contribution points. For the first contribution point, as an example, we recall a well-known adaptive technique for the learning rate  $\eta_m$ :  $\eta_m = (1 + a)\eta_{m-1}$  if the error is decreasing, and  $\eta_m = (1 - a)\eta_{m-1}$  otherwise, where  $a < 1$  is a positive number. Xu's condition  $\lim_{m \rightarrow \infty} \eta_m / \eta_{m+1} = 1$  (Xu et al., 2009) is not valid in this case, while our convergence results remains valid. For the second contribution point, it is interesting to see that the learning method OGM-SS with stochastic nature enjoys deterministic convergence. We observe that OGM-F is actually a deterministic iteration procedure in that the iteration sequence is determined uniquely by the initial value and the fixed order of the samples. The convergence result is a bit easier to prove for OGM-F than for OGM-SS. We have reason to believe, and our experience shows, that OGM-SS behaves numerically better than OGM-F since the stochastic nature of the learning procedure survives in OGM-SS (Li & Ding, 2005; Li et al., 2004). Our convergence results are generalizations of both the results of

Xu et al. (2009), which considers OGM-F, and the results of Li et al. (2004), which considers OGM-SS with an unpleasant condition  $\eta_m = O(1/m)$  on the learning rate. Our third contribution allows the activation functions for both hidden and output layers to be more flexible. Here we remark that typically, S–S type networks are used for classification problems, and S–P type networks with Sigmoid hidden neurons and linear output neurons are used for approximation problems. The existing convergence results (Li et al., 2004; Shao et al., 2007; Wu & Xu, 2002; Wu et al., 2005, 2002; Wu & Shao, 2003; Wu et al., 2005; Xu et al., 2009) are mostly for either S–S type or S–P type alone but not for both of them. In this paper, we give a uniform treatment for all types of networks. The fourth and fifth contribution points are mainly of theoretical interest. From a theoretical point of view, we mention that different analytical tools are employed in Wu et al. (2005) and Xu et al. (2009) and this study for the convergence analysis, might explain, at least in part, why different conditions are obtained for the convergence. The differential Taylor expansion is used in Wu et al. (2005), which requires the boundedness of the second derivative  $g''$  of the activation function  $g$ ; the mean value theorem of integrals is employed in Xu et al. (2009), which requires  $g'$  to be Lipschitz continuous and uniformly bounded; and in this paper, we use the integral Taylor expansion and hence require the Lipschitz continuity of  $g'$  on any bounded closed interval. Finally, we point out that Xu et al. (2009) is a big step forward for the convergence study of OGM-F and that Xu et al. (2009) also includes another convergence result under the condition that the error function is directionally convex. This convex condition is not considered in this paper.

The rest of this paper is organized as follows. In Section 2, online updating methods including OGM-F and OGM-SS are introduced. The main convergence results are presented in Section 3 and their proofs are gathered in Section 4. Some conclusions are drawn in Section 5.

## 2. OGM-F and OGM-SS

Let us begin with an introduction of a feedforward neural network with three layers. The numbers of neurons for the input, hidden and output layers are  $p$ ,  $n$  and 1, respectively. Suppose that the training sample set is  $\{\mathbf{x}^j, O^j\}_{j=1}^J \subset \mathbb{R}^p \times \mathbb{R}$ , where  $\mathbf{x}^j$  and  $O^j$  are the input and the corresponding ideal output of the  $j$ th sample, respectively. Let  $\mathbf{V} = (v_{i,j})_{n \times p}$  be the weight matrix connecting the input and the hidden layers, and write  $\mathbf{v}_i = (v_{i1}, v_{i2}, \dots, v_{ip})^T$  for  $i = 1, 2, \dots, n$ . The weight vector connecting the hidden and the output layers is denoted by  $\mathbf{u} = (u_1, u_2, \dots, u_n)^T \in \mathbb{R}^n$ . To simplify the presentation, we combine the weight matrix  $\mathbf{V}$  with the weight vector  $\mathbf{u}$ , and write  $\mathbf{w} = (\mathbf{u}^T, \mathbf{v}_1^T, \dots, \mathbf{v}_n^T)^T \in \mathbb{R}^{n(p+1)}$ . Let  $g, f : \mathbb{R} \rightarrow \mathbb{R}$  be given activation functions for the hidden and output layers, respectively. For convenience, we introduce the following vector valued function

$$G(\mathbf{z}) = (g(z_1), g(z_2), \dots, g(z_n))^T, \quad \forall \mathbf{z} \in \mathbb{R}^n. \quad (1)$$

For any given input  $\mathbf{x} \in \mathbb{R}^p$ , the output of the hidden neurons is  $G(\mathbf{V}\mathbf{x})$ , and the final actual output is

$$y = f(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x})). \quad (2)$$

For any fixed weights  $\mathbf{w}$ , the error of the neural networks is defined as

$$E(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^J (O^j - f(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)))^2 = \sum_{j=1}^J f_j(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)), \quad (3)$$

where  $f_j(t) = \frac{1}{2}(O^j - f(t))^2$ ,  $j = 1, 2, \dots, J$ ,  $t \in \mathbb{R}$ . The gradients of the error function with respect to  $\mathbf{u}$  and  $\mathbf{v}_i$  are, respectively,

given by

$$\begin{aligned} E_{\mathbf{u}}(\mathbf{w}) &= -\sum_{j=1}^J (O^j - y^j) f'(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)) G(\mathbf{V}\mathbf{x}^j) \\ &= \sum_{j=1}^J f'_j(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)) G(\mathbf{V}\mathbf{x}^j), \end{aligned} \quad (4)$$

$$\begin{aligned} E_{\mathbf{v}_i}(\mathbf{w}) &= -\sum_{j=1}^J (O^j - y^j) f'(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)) u_i g'(\mathbf{v}_i \cdot \mathbf{x}^j) \mathbf{x}^j \\ &= \sum_{j=1}^J f'_j(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)) u_i g'(\mathbf{v}_i \cdot \mathbf{x}^j) \mathbf{x}^j, \end{aligned} \quad (5)$$

where

$$y^j = f(\mathbf{u} \cdot G(\mathbf{V}\mathbf{x}^j)), \quad i = 1, 2, \dots, n; j = 1, 2, \dots, J. \quad (6)$$

Write

$$E_{\mathbf{v}}(\mathbf{w}) = (E_{\mathbf{v}_1}(\mathbf{w})^T, E_{\mathbf{v}_2}(\mathbf{w})^T, \dots, E_{\mathbf{v}_n}(\mathbf{w})^T)^T, \quad (7)$$

$$E_{\mathbf{w}}(\mathbf{w}) = (E_{\mathbf{u}}(\mathbf{w})^T, E_{\mathbf{v}}(\mathbf{w})^T)^T. \quad (8)$$

First, let us consider the case that the training samples are supplied to the network in a fixed order (OGM-F) in the training process. Hence, starting from an arbitrary initial guess  $\mathbf{w}^0$ , we proceed to refine it iteratively by the formulas

$$\mathbf{u}^{m+j+1} = \mathbf{u}^{m+j} + \Delta_j \mathbf{u}^{m+j}, \quad (9)$$

$$\mathbf{v}_i^{m+j+1} = \mathbf{v}_i^{m+j} + \Delta_j \mathbf{v}_i^{m+j}, \quad (10)$$

where

$$\begin{aligned} \Delta_k \mathbf{u}^{m+j} &= \eta_m (O^k - y^{m+j, k}) f'(\mathbf{u}^{m+j} \cdot G^{m+j, k}) G^{m+j, k} \\ &= -\eta_m f'_k(\mathbf{u}^{m+j} \cdot G^{m+j, k}) G^{m+j, k}, \end{aligned} \quad (11)$$

$$\begin{aligned} \Delta_k \mathbf{v}_i^{m+j} &= \eta_m (O^k - y^{m+j, k}) f'(\mathbf{u}^{m+j} \cdot G^{m+j, k}) \\ &\quad \times u_i^{m+j} g'(\mathbf{v}_i^{m+j} \cdot \mathbf{x}^k) \mathbf{x}^k \\ &= -\eta_m f'_k(\mathbf{u}^{m+j} \cdot G^{m+j, k}) u_i^{m+j} g'(\mathbf{v}_i^{m+j} \cdot \mathbf{x}^k) \mathbf{x}^k, \end{aligned} \quad (12)$$

$$G^{m+j, k} = G(\mathbf{V}^{m+j} \mathbf{x}^k), \quad y^{m+j, k} = f(\mathbf{u}^{m+j} \cdot G^{m+j, k}), \quad (13)$$

Here the parameter  $\eta_m$  is the learning rate, whose value may be changed after each cycle of the training procedure.

We can also choose training samples in a special stochastic order (OGM-SS) as follows: For the  $m$ th training cycle, let  $\{\mathbf{x}^{m,1}, \mathbf{x}^{m,2}, \dots, \mathbf{x}^{m,J}\}$  be a stochastic permutation of the set  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^J\}$ . Similar to (9) and (10), the weights are iteratively updated in the following fashion

$$\mathbf{u}^{m+j+1} = \mathbf{u}^{m+j} + \Delta_j^m \mathbf{u}^{m+j}, \quad (14)$$

$$\mathbf{v}_i^{m+j+1} = \mathbf{v}_i^{m+j} + \Delta_j^m \mathbf{v}_i^{m+j}, \quad (15)$$

where

$$\begin{aligned} \Delta_k^m \mathbf{u}^{m+j} &= \eta_m (O^k - y^{m+j, m, k}) f'(\mathbf{u}^{m+j} \cdot G^{m+j, m, k}) G^{m+j, m, k} \\ &= -\eta_m f'_k(\mathbf{u}^{m+j} \cdot G^{m+j, m, k}) G^{m+j, m, k}, \end{aligned} \quad (16)$$

$$\begin{aligned} \Delta_k^m \mathbf{v}_i^{m+j} &= \eta_m (O^k - y^{m+j, m, k}) f'(\mathbf{u}^{m+j} \cdot G^{m+j, m, k}) \cdot u_i^{m+j} \\ &\quad \times g'(\mathbf{v}_i^{m+j} \cdot \mathbf{x}^{m, k}) \mathbf{x}^{m, k} \\ &= -\eta_m f'_k(\mathbf{u}^{m+j} \cdot G^{m+j, m, k}) \\ &\quad \times u_i^{m+j} g'(\mathbf{v}_i^{m+j} \cdot \mathbf{x}^{m, k}) \mathbf{x}^{m, k}, \end{aligned} \quad (17)$$

$$G^{m+j, m, k} = G(\mathbf{V}^{m+j} \mathbf{x}^{m, k}), \quad y^{m+j, m, k} = f(\mathbf{u}^{m+j} \cdot G^{m+j, m, k}), \quad (18)$$

We mention that OGM-F and OGM-SS are also called cycle learning and almost-cycle learning in Heskes and Wiegerinck (1996), respectively.

### 3. Main results

For any  $\mathbf{x} \in \mathbb{R}^n$ , we write  $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^n x_i^2}$ , where  $\|\cdot\|$  stands for the Euclidean norm in  $\mathbb{R}^n$ . Let  $\Omega_0 = \{\mathbf{w} \in \Omega : E_{\mathbf{w}}(\mathbf{w}) = 0\}$  be the stationary point set of the error function  $E(\mathbf{w})$ , where  $\Omega \subset \mathbb{R}^{n(p+1)}$  is a bounded region satisfying (A3) below. Let  $\Omega_{0,s} \subset \mathbb{R}$  be the projection of  $\Omega_0$  onto the  $s$ th coordinate axis, that is,

$$\Omega_{0,s} = \{w_s \in \mathbb{R} : \mathbf{w} = (w_1, \dots, w_s, \dots, w_{n(p+1)})^T \in \Omega_0\} \quad (19)$$

for  $s = 1, 2, \dots, n(p+1)$ . To analyze the convergence of the algorithm, we need the following assumptions.

- (A1)  $g'(t)$  and  $f'(t)$  are Lipschitz continuous on any bounded closed interval;
- (A2)  $\eta_m > 0$ ,  $\sum_{m=0}^{\infty} \eta_m = \infty$ ,  $\sum_{m=0}^{\infty} \eta_m^2 < \infty$ ;
- (A3) There exists a bounded open set  $\Omega \subset \mathbb{R}^n$  such that  $\{\mathbf{w}^m\} \subset \Omega$  ( $m \in \mathbb{N}$ );
- (A3') There exists a bounded open set  $\Omega' \subset \mathbb{R}^n$  such that  $\{\mathbf{u}^m\} \subset \Omega'$  ( $m \in \mathbb{N}$ ), and the derivative of the activation function  $g$  in (1) is uniformly bounded and Lipschitz continuous on  $\mathbb{R}$ .
- (A4)  $\Omega_{0,s}$  does not contain any interior point for every  $s = 1, 2, \dots, n(p+1)$ .

**Theorem 3.1.** Assume that conditions (A1)–(A3) are valid. Then, starting from an arbitrary initial value  $\mathbf{w}^0$ , the learning sequence  $\{\mathbf{w}^m\}$  defined by (9) and (10) or by (14) and (15) satisfies the following weak convergence

$$\lim_{m \rightarrow \infty} \|E_{\mathbf{w}}(\mathbf{w}^m)\| = 0. \quad (20)$$

Moreover, if assumptions (A1)–(A4) are valid, there holds the strong convergence: There exists  $\mathbf{w}^* \in \Omega_0$  such that

$$\lim_{m \rightarrow \infty} \mathbf{w}^m = \mathbf{w}^*. \quad (21)$$

Let us make three remarks on the convergence result: (1) We claim that the weak convergence remains valid if the activation function  $g$  of the hidden layer is a commonly used sigmoid function and assumptions (A3') (instead of (A3)) and (A2) are valid. This is due to the fact that the sigmoid function  $g$  is uniformly bounded on  $\mathbb{R}$  and that (37) is valid even if the weight vectors  $\mathbf{v}_i$  ( $i = 1, 2, \dots, n$ ) are unbounded. (2) In the numerical analysis of an iterative method for a class of nonlinear problems, the iterative sequence is often required to be bounded in order to prove its convergence. This is what we do in conditions (A3) and (A3'). We mention that the weights will be automatically bounded in the network training with the help of a penalty term (cf. Zhang et al., 2009). (3) For the strong convergence, our condition (A4) on  $\Omega_0$  allows it to be finite set, countably infinite set, nowhere dense set or even some uncountable dense set. Hence, the corresponding assumptions that the set  $\Omega_0$  contains finite points and at most countably infinite points in Wu et al. (2005) and Xu et al. (2009), respectively, are special cases of assumption (A4). This relaxed condition makes it much easier to verify the strong convergence in practice.

### 4. Proofs

For convenience of presentation, we present in detail the convergence proof for OGM-F in the following Section 4.1. Then, in Section 4.2, we briefly point out how to extend the result to OGM-SS.

#### 4.1. Convergence analysis for OGM-F

We first present four useful lemmas for the convergence analysis.



**Lemma 4.1.** Let  $q(x)$  be a function defined on a bounded closed interval  $[a, b]$  such that  $q'(x)$  is Lipschitz continuous with Lipschitz constant  $K > 0$ . Then,  $q'(x)$  is differentiable almost everywhere in  $[a, b]$  and

$$|q''(x)| \leq K, \quad \text{a.e. } [a, b]. \quad (22)$$

Moreover, there exists a constant  $C > 0$  such that

$$q(x) \leq q(x_0) + q'(x_0)(x - x_0) + C(x - x_0)^2, \quad \forall x_0, x \in [a, b]. \quad (23)$$

**Proof.** Since  $q'(x)$  is Lipschitz continuous on  $[a, b]$ ,  $q'(x)$  is absolutely continuous and the derivative  $q''(x)$  exists almost everywhere on  $[a, b]$ . Hence, for almost every  $x \in [a, b]$ ,

$$\begin{aligned} |q''(x)| &= \left| \lim_{h \rightarrow 0} \frac{q'(x+h) - q'(x)}{h} \right| \\ &= \lim_{h \rightarrow 0} \left| \frac{q'(x+h) - q'(x)}{h} \right| \leq K. \end{aligned} \quad (24)$$

Using the integral Taylor expansion, we deduce that

$$\begin{aligned} q(x) &= q(x_0) + q'(x_0)(x - x_0) \\ &\quad + (x - x_0)^2 \int_0^1 (1-t)q''(x_0 + t(x - x_0))dt \\ &\leq q(x_0) + q'(x_0)(x - x_0) + (x - x_0)^2 \int_0^1 K(1-t)dt \\ &= q(x_0) + q'(x_0)(x - x_0) + C(x - x_0)^2, \\ &\quad \left( C = \frac{K}{2}, x_0, x \in [a, b] \right). \quad \square \end{aligned} \quad (25)$$

**Lemma 4.2.** Suppose that the learning rate  $\eta_m$  satisfies (A2) and that the sequence  $\{a_m\}$  ( $m \in \mathbb{N}$ ) satisfies  $a_m \geq 0$ ,  $\sum_{m=0}^{\infty} \eta_m a_m^\beta < \infty$  and  $|a_{m+1} - a_m| \leq \mu \eta_m$  for some positive constants  $\beta$  and  $\mu$ . Then we have

$$\lim_{m \rightarrow \infty} a_m = 0. \quad (26)$$

**Proof.** According to (A2), we know that  $\eta_m \rightarrow 0$  as  $m \rightarrow \infty$ . We claim that  $\lim_{k \rightarrow \infty} \inf_{m > k} a_m = 0$ . Otherwise, if  $a_* \equiv \lim_{k \rightarrow \infty} \inf_{m > k} a_m \in (0, \infty]$ , then by the definition of the inferior limit, there exists an integer  $M > k$  such that  $a_m \geq \frac{a_*}{2} > 0$  for  $m \geq M$ , which leads to

$$\sum_{m=0}^{\infty} \eta_m a_m^\beta \geq \left(\frac{a_*}{2}\right)^\beta \sum_{m=M}^{\infty} \eta_m = \infty. \quad (27)$$

This contradicts  $\sum_{m=0}^{\infty} \eta_m a_m^\beta < \infty$  and confirms the claim. Next, we claim that  $\lim_{k \rightarrow \infty} \sup_{m > k} a_m = 0$ . Otherwise, there exists  $\delta \in (0, \infty]$  such that  $\lim_{k \rightarrow \infty} \sup_{m > k} a_m = \delta$ . Then, for any  $0 < \varepsilon < \delta$ , we can choose two subsequences  $\{a_{i_k}\}$  and  $\{a_{j_k}\}$  of  $\{a_m\}$  to satisfy (1)  $a_{i_k} \in (0, \frac{\varepsilon}{4})$ ,  $a_{j_k} \in (\varepsilon, \delta)$ ; (2)  $i_k + 1 < j_k < i_{k+1}$ ; (3)  $a_{i_k+1} \in [\frac{\varepsilon}{4}, \frac{\varepsilon}{2}]$ . (This can be done because  $\lim_{k \rightarrow \infty} \inf_{m > k} a_m = 0$ ,  $\lim_{k \rightarrow \infty} \sup_{m > k} a_m = \delta$ , and  $|a_m - a_{m+1}| \leq \mu \eta_m \rightarrow 0$  as  $m \rightarrow \infty$ .) For any  $i_k < m < j_k$ , we have  $a_m \in [\frac{\varepsilon}{4}, \varepsilon]$ . Thus, we conclude that

$$\begin{aligned} \frac{\varepsilon}{2} &\leq |a_{j_k} - a_{i_k+1}| \leq |a_{j_k} - a_{j_k-1}| + \dots + |a_{i_k+2} - a_{i_k+1}| \\ &\leq \mu \sum_{m=i_k+1}^{j_k-1} \eta_m \leq \mu \sum_{m=i_k+1}^{j_k} \eta_m. \end{aligned}$$

Therefore, we have for all large enough integer  $k$  that

$$\sum_{m=i_k}^{j_k} \eta_m a_m^\beta \geq \sum_{m=i_k+1}^{j_k} \eta_m a_m^\beta \geq \left(\frac{\varepsilon}{4}\right)^\beta \sum_{m=i_k+1}^{j_k} \eta_m \geq \frac{2}{\mu} \left(\frac{\varepsilon}{4}\right)^{\beta+1}.$$

But this contradicts  $\sum_{m=0}^{\infty} \eta_m a_m^\beta < \infty$  and implies our second claim. Finally, the above two claims together clearly lead to the desired estimate (26).  $\square$

**Lemma 4.3.** Let  $\{b_m\}$  be a bounded sequence satisfying  $\lim_{m \rightarrow \infty} (b_{m+1} - b_m) = 0$ . Write  $\gamma_1 = \lim_{n \rightarrow \infty} \inf_{m > n} b_m$ ,  $\gamma_2 = \lim_{n \rightarrow \infty} \sup_{m > n} b_m$  and  $S = \{a \in \mathbb{R} : \text{There exists a subsequence } \{b_{i_k}\} \text{ of } \{b_m\} \text{ such that } b_{i_k} \rightarrow a \text{ as } k \rightarrow \infty\}$ . Then we have

$$S = [\gamma_1, \gamma_2]. \quad (28)$$

**Proof.** It is obvious that  $\gamma_1 \leq \gamma_2$  and  $S \subseteq [\gamma_1, \gamma_2]$ . If  $\gamma_1 = \gamma_2$ , then (28) follows simply from  $\lim_{m \rightarrow \infty} b_m = \gamma_1 = \gamma_2$ . Let us consider the case  $\gamma_1 < \gamma_2$  and proceed to prove that  $S \supseteq [\gamma_1, \gamma_2]$ .

For any  $a \in (\gamma_1, \gamma_2)$ , there exists  $\varepsilon > 0$  such that  $(a - \varepsilon, a + \varepsilon) \subseteq (\gamma_1, \gamma_2)$ . Noting that  $\lim_{m \rightarrow \infty} (b_{m+1} - b_m) = 0$ , we observe that  $b_m$  travels to and from between  $\gamma_1$  and  $\gamma_2$  with very small pace for all large enough  $m$ . Hence, there must be infinite number of points of the sequence  $\{b_m\}$  falling into  $(a - \varepsilon, a + \varepsilon)$ . This implies  $a \in S$  and thus  $(\gamma_1, \gamma_2) \subseteq S$ . Furthermore,  $(\gamma_1, \gamma_2) \subseteq S$  immediately leads to  $[\gamma_1, \gamma_2] \subseteq S$ . This completes the proof.  $\square$

Let the sequence  $\{\mathbf{w}^{m+j}\}$  ( $m \in \mathbb{N}$ ,  $j = 1, 2, \dots, J$ ) be generated by (9) and (10). We introduce the following notations:

$$\mathbf{R}^{m,j} = \Delta_j \mathbf{u}^{m+j} - \Delta_j \mathbf{u}^{m,j}, \quad (29)$$

$$\mathbf{r}_i^{m,j} = \Delta_j \mathbf{v}_i^{m+j} - \Delta_j \mathbf{v}_i^{m,j}, \quad (30)$$

$$\mathbf{d}^{m,l} = \mathbf{u}^{m+l} - \mathbf{u}^{m,j} = \sum_{j=1}^l \Delta_j \mathbf{u}^{m+j} = \sum_{j=1}^l \Delta_j \mathbf{u}^{m,j} + \sum_{j=1}^l \mathbf{R}^{m,j}, \quad (31)$$

$$\mathbf{h}_i^{m,l} = \mathbf{v}_i^{m+l} - \mathbf{v}_i^{m,j} = \sum_{j=1}^l \Delta_j \mathbf{v}_i^{m+j} = \sum_{j=1}^l \Delta_j \mathbf{v}_i^{m,j} + \sum_{j=1}^l \mathbf{r}_i^{m,j}, \quad (32)$$

$$\psi^{m,l,j} = \mathbf{G}^{m+l,j} - \mathbf{G}^{m,j}, \quad (33)$$

$$m \in \mathbb{N}, j = 1, 2, \dots, J, l = 1, 2, \dots, J, i = 1, 2, \dots, n.$$

Then, (9) and (10) can be rewritten as

$$\mathbf{u}^{m+j} = \mathbf{u}^{m,j} + \sum_{k=1}^j (\Delta_k \mathbf{u}^{m,k} + \mathbf{R}^{m,k}), \quad (34)$$

$$\mathbf{v}_i^{m+j} = \mathbf{v}_i^{m,j} + \sum_{k=1}^j (\Delta_k \mathbf{v}_i^{m,k} + \mathbf{r}_i^{m,k}). \quad (35)$$

Let constants  $C_1$  and  $C_2$  be defined by (cf. assumption (A3))

$$\max_{1 \leq j \leq J} \{\|\mathbf{x}^j\|, |\mathcal{O}^j|\} = C_1, \quad \sup_{m \in \mathbb{N}} \|\mathbf{w}^m\| = C_2. \quad (36)$$

By assumption (A1),  $f'_j(t)$  also satisfies the Lipschitz condition for  $j = 1, 2, \dots, J$ . Furthermore,  $g(t)$ ,  $f(t)$  and  $f_j(t)$  are all uniformly continuous on any bounded closed interval.

**Lemma 4.4.** Let conditions (A1) and (A3) be valid, and let the sequence  $\{\mathbf{w}^{m+j}\}$  be generated by (9) and (10). Then there are constants  $C_3 - C_7$  such that

$$\|\mathbf{G}^{m+j,k}\| \leq C_3, \quad (37)$$

$$\|\mathbf{d}^{m,l}\| \leq C_4 \eta_m, \quad \|\psi^{m,l,j}\| \leq C_5 \eta_m, \quad (38)$$

$$\|\mathbf{R}^{m,j}\| \leq C_6 \eta_m^2, \quad \|\mathbf{r}_i^{m,j}\| \leq C_7 \eta_m^2, \quad (39)$$

where  $m \in \mathbb{N}$ ;  $j, k = 1, 2, \dots, J$ ;  $l = 1, 2, \dots, J$ ;  $i = 1, 2, \dots, n$ .

**Proof.** According to (36), we have

$$|\mathbf{v}_i^{mj+j} \cdot \mathbf{x}^k| \leq \|\mathbf{v}_i^{mj+j}\| \|\mathbf{x}^k\| \leq C_1 C_2 \equiv D_1. \quad (40)$$

Thus, there exists a positive constant  $C_{3,1}$  such that

$$\max_{|t| \leq D_1} |g(t)| = C_{3,1}, \quad (41)$$

$$\|G^{mj+j, k}\| = \|G(\mathbf{v}^{mj+j} \mathbf{x}^k)\| \leq \sqrt{n} C_{3,1} \equiv C_3. \quad (42)$$

It follows from (36) and (42) that

$$|\mathbf{u}^{mj+j} \cdot G^{mj+j, k}| \leq \|\mathbf{u}^{mj+j}\| \|G^{mj+j, k}\| \leq C_2 C_3 \equiv D_2. \quad (43)$$

Then, there is a positive constant  $C_{4,1}$  such that

$$\max_{|t| \leq D_2} |f_j'(t)| \leq C_{4,1}. \quad (44)$$

Furthermore, a combination of (A1), (11), (37) and (40) gives

$$\|d^{m, l}\| = \|\mathbf{u}^{m+l} - \mathbf{u}^{mj}\| = \left\| \sum_{j=1}^l \Delta_j \mathbf{u}^{mj+j} \right\| \leq C_4 \eta_m, \quad (45)$$

where  $C_4 = J C_{4,1} C_3$ .

Employing (40), we find that

$$\max_{|t| \leq D_1} |g'(t)| = C_{5,1}, \quad (46)$$

where  $C_{5,1}$  is a positive constant. Moreover, we observe that

$$\begin{aligned} \|\psi^{m, l, j}\| &= \|G^{m+l, j} - G^{mj, j}\| \leq \max_{1 \leq i \leq n} |g'(t_i)| \|\mathbf{x}^j\| \sum_{i=1}^n \|\mathbf{h}_i^{m, l}\| \\ &\leq \max_{1 \leq i \leq n} |g'(t_i)| \|\mathbf{x}^j\| \sum_{i=1}^n \sum_{k=1}^l \|\Delta_k \mathbf{v}_i^{mj+k}\| \\ &\leq C_5 \eta_m, \end{aligned} \quad (47)$$

where  $C_5 = n l C_{4,1} C_{5,1} \max_{1 \leq i \leq n} |g'(t_i)| \|\mathbf{x}^j\| \sup_{m \in \mathbb{N}} \|\mathbf{w}^m\| \max_{1 \leq k \leq l} \|\mathbf{x}^k\|$ , in which  $t_i = \mathbf{v}_i^{mj} \cdot \mathbf{x}^j + \theta_i (\mathbf{v}_i^{mj+l} - \mathbf{v}_i^{mj}) \cdot \mathbf{x}^j$ ,  $\theta_i \in (0, 1)$ , and  $|t_i| \leq |\mathbf{v}_i^{mj} \cdot \mathbf{x}^j| + |(\mathbf{v}_i^{mj+l} - \mathbf{v}_i^{mj}) \cdot \mathbf{x}^j| \leq 3 C_1 C_2$ . By virtue of (A1), we see that  $|g'(t_i)|$  ( $i = 1, 2, \dots, n$ ) is bounded.

Combining  $f_j'(t)$ 's Lipschitz continuity, (36) and (37), we have

$$\begin{aligned} &|f_j'(\mathbf{u}^{mj+j} \cdot G^{mj+j, j}) - f_j'(\mathbf{u}^{mj} \cdot G^{mj, j})| \\ &\leq L |\mathbf{u}^{mj+j} \cdot G^{mj+j, j} - \mathbf{u}^{mj} \cdot G^{mj, j}| \\ &\leq L \|d^{m, j}\| \|G^{mj+j, j}\| \leq L C_3 \|d^{m, j}\|, \end{aligned} \quad (48)$$

$$\begin{aligned} &|f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) - f_j'(\mathbf{u}^{mj} \cdot G^{mj, j})| \\ &\leq L |\mathbf{u}^{mj} \cdot G^{mj+j, j} - \mathbf{u}^{mj} \cdot G^{mj, j}| \\ &\leq L \|\mathbf{u}^{mj}\| \|\psi^{m, j, j}\| \leq L C_2 \|\psi^{m, j, j}\|, \end{aligned} \quad (49)$$

where  $L > 0$  is the Lipschitz constant.

By the definition of  $R^{m, j}$ , we see that

$$\begin{aligned} R^{m, j} &= \Delta_j \mathbf{u}^{mj+j} - \Delta_j \mathbf{u}^{mj} \\ &= -\eta_m (f_j'(\mathbf{u}^{mj+j} \cdot G^{mj+j, j}) G^{mj+j, j} - f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) G^{mj, j}) \\ &= -\eta_m [f_j'(\mathbf{u}^{mj+j} \cdot G^{mj+j, j}) \psi^{m, j, j} \\ &\quad + (f_j'(\mathbf{u}^{mj+j} \cdot G^{mj+j, j}) - f_j'(\mathbf{u}^{mj} \cdot G^{mj, j})) G^{mj, j} \\ &\quad + (f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) - f_j'(\mathbf{u}^{mj} \cdot G^{mj, j})) G^{mj, j}]. \end{aligned} \quad (50)$$

Therefore, it follows from (37), (38), (48) and (49) that

$$\|R^{m, j}\| \leq \eta_m (L C_3^2 \|d^{m, j}\| + (C_{4,1} + L C_2 C_3) \|\psi^{m, j, j}\|) \leq C_6 \eta_m^2, \quad (51)$$

where  $C_6 = \max\{L C_3^2 C_4, (C_{4,1} + L D_2) C_5\}$ .

Similarly, we can show the existence of a constant  $C_7 > 0$  such that

$$\|r_i^{m, j}\| \leq C_7 \eta_m^2. \quad \square \quad (52)$$

The next lemma reveals an almost monotonicity of the error function during the training process.

**Lemma 4.5.** Let the sequence  $\{\mathbf{w}^{m+j}\}$  be generated by (9) and (10). Under assumptions (A1) and (A3), there holds

$$E(\mathbf{w}^{(m+1)J}) \leq E(\mathbf{w}^{mj}) - \eta_m \|E_{\mathbf{w}}(\mathbf{w}^{mj})\|^2 + C_8 \eta_m^2, \quad (53)$$

$$(m = 0, 1, \dots)$$

where  $C_8 > 0$  is a constant independent of  $m$  and  $\eta_m$ .

**Proof.** By virtue of assumption (A1) and Lemma 4.1, we know that the derivative  $g''(\mathbf{v}_i^{mj} \cdot \mathbf{x}^j + t(\mathbf{h}_i^{m, J} \cdot \mathbf{x}^j))$  is integrable almost everywhere on  $[0, 1]$  and

$$\begin{aligned} &f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) \mathbf{u}^{mj} \cdot \psi^{m, J, j} \\ &= f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) \sum_{i=1}^n u_i^{mj} g'(\mathbf{v}_i^{mj} \cdot \mathbf{x}^j) h_i^{m, J} \cdot \mathbf{x}^j \\ &\quad + f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) \sum_{i=1}^n u_i^{mj} (h_i^{m, J} \cdot \mathbf{x}^j)^2 \cdot \int_0^1 (1-t) \\ &\quad \times g''(\mathbf{v}_i^{mj} \cdot \mathbf{x}^j + t(\mathbf{h}_i^{m, J} \cdot \mathbf{x}^j)) dt. \end{aligned} \quad (54)$$

By virtue of Lemma 4.1, (11), (12) and (54), there is a constant  $C_9 > 0$  such that

$$\begin{aligned} &f_j(\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J, j}) \\ &\leq f_j(\mathbf{u}^{mj} \cdot G^{mj, j}) \\ &\quad + f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) (\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J, j} - \mathbf{u}^{mj} \cdot G^{mj, j}) \\ &\quad + C_9 (\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J, j} - \mathbf{u}^{mj} \cdot G^{mj, j})^2 \\ &= f_j(\mathbf{u}^{mj} \cdot G^{mj, j}) \\ &\quad + f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) (d^{m, J} \cdot G^{mj, j} + \mathbf{u}^{mj} \cdot \psi^{m, J, j} + d^{m, J} \cdot \psi^{m, J, j}) \\ &\quad + C_9 (\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J, j} - \mathbf{u}^{mj} \cdot G^{mj, j})^2 \\ &= f_j(\mathbf{u}^{mj} \cdot G^{mj, j}) - \frac{1}{\eta_m} \Delta_j \mathbf{u}^{mj} \cdot d^{m, J} - \frac{1}{\eta_m} \sum_{i=1}^n (\Delta_j \mathbf{v}_i^{mj} \cdot h_i^{m, J}) \\ &\quad + f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) \sum_{i=1}^n u_i^{mj} (h_i^{m, J} \cdot \mathbf{x}^j)^2 \cdot \int_0^1 (1-t) \\ &\quad \times g''(\mathbf{v}_i^{mj} \cdot \mathbf{x}^j + t(\mathbf{h}_i^{m, J} \cdot \mathbf{x}^j)) dt \\ &\quad + f_j'(\mathbf{u}^{mj} \cdot G^{mj, j}) d^{m, J} \cdot \psi^{m, J, j} \\ &\quad + C_9 (\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J, j} - \mathbf{u}^{mj} \cdot G^{mj, j})^2. \end{aligned} \quad (55)$$

Summing (55) from  $j = 1$  to  $j = J$  and noting (3)–(5), (31) and (32), we have

$$\begin{aligned} E(\mathbf{w}^{(m+1)J}) &\leq E(\mathbf{w}^{mj}) - \frac{1}{\eta_m} \left( \left\| \sum_{j=1}^J \Delta_j \mathbf{u}^{mj} \right\|^2 + \sum_{i=1}^n \left\| \sum_{j=1}^J \Delta_j \mathbf{v}_i^{mj} \right\|^2 \right) + \delta_m \\ &= E(\mathbf{w}^{mj}) - \eta_m \left( \|E_{\mathbf{u}}(\mathbf{w}^{mj})\|^2 + \sum_{i=1}^n \|E_{\mathbf{v}_i}(\mathbf{w}^{mj})\|^2 \right) + \delta_m \\ &= E(\mathbf{w}^{mj}) - \eta_m \|E_{\mathbf{w}}(\mathbf{w}^{mj})\|^2 + \delta_m, \end{aligned} \quad (56)$$

where

$$\begin{aligned} \delta_m = & -\frac{1}{\eta_m} \sum_{j=1}^J \Delta_j \mathbf{u}^{mj} \cdot \sum_{j=1}^J R^{m,j} - \frac{1}{\eta_m} \sum_{i=1}^n \left( \sum_{j=1}^J \Delta_j \mathbf{v}_i^{mj} \cdot \sum_{j=1}^J r_i^{m,j} \right) \\ & + \sum_{j=1}^J \sum_{i=1}^n u_i^{mj} f'_j(\mathbf{u}^{mj} \cdot G^{mj,j}) (h_i^{m,J} \cdot \mathbf{x}^j)^2 \cdot \int_0^1 (1-t) \\ & \times g''(\mathbf{v}_i^{mj} \cdot \mathbf{x}^j + t(h_i^{m,J} \cdot \mathbf{x}^j)) dt \\ & + \sum_{j=1}^J f'_j(\mathbf{u}^{mj} \cdot G^{mj,j}) d^{m,J} \cdot \psi^{m,J,j} \\ & + C_9 \sum_{j=1}^J (\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j} - \mathbf{u}^{mj} \cdot G^{mj,j})^2. \end{aligned}$$

It now follows from (36) and (37) that

$$\begin{aligned} \|G^{mj,j}\| = \|G(\mathbf{V}^{mj} \mathbf{x}^j)\| & \leq C_3, \\ |\mathbf{u}^{mj} \cdot G^{mj,j}| & \leq \|\mathbf{u}^{mj}\| \|G^{mj,j}\| \leq C_2 C_3 = D_2. \end{aligned} \quad (57)$$

By (11), (42)–(44) and (51), the first term of  $\delta_m$  can be estimated as follows.

$$\begin{aligned} & \left\| -\frac{1}{\eta_m} \sum_{j=1}^J \Delta_j \mathbf{u}^{mj} \cdot \sum_{j=1}^J R^{m,j} \right\| \\ & \leq \frac{1}{\eta_m} \sum_{j=1}^J \|\Delta_j \mathbf{u}^{mj}\| \cdot \sum_{j=1}^J \|R^{m,j}\| \leq C_{8,1} \eta_m^2, \end{aligned} \quad (58)$$

where  $C_{8,1} = J^2 C_3 C_{4,1} C_6 = J C_4 C_6$ .

Similar estimates for the other terms of  $\delta_m$  can be obtained with corresponding constants  $C_{8,t} > 0$  for  $t = 2, \dots, 5$ . Finally, the desired estimate (53) is proved by setting  $C_8 = \sum_{t=1}^5 C_{8,t}$ .  $\square$

Now, we are ready to prove the convergence theorem.

**Proof of Theorem 3.1 for OGM-F.** The proof is divided into two parts, dealing with (20) and (21), respectively.

**Proof of (20).** By (A2) and Lemma 4.5, we conclude that

$$\begin{aligned} & \sum_{m=0}^{\infty} \eta_m \|E_{\mathbf{w}}(\mathbf{w}^{mj})\|^2 \\ & = \sum_{m=0}^{\infty} \eta_m \left( \|E_{\mathbf{u}}(\mathbf{w}^{mj})\|^2 + \|E_{\mathbf{v}}(\mathbf{w}^{mj})\|^2 \right) < \infty, \end{aligned} \quad (59)$$

$$\sum_{m=0}^{\infty} \eta_m \|E_{\mathbf{u}}(\mathbf{w}^{mj})\|^2 < \infty. \quad (60)$$

Employing the integral Taylor expansion, we deduce that

$$\begin{aligned} & f'_j(\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j}) G^{(m+1)J,j} - f'_j(\mathbf{u}^{mj} \cdot G^{mj,j}) G^{mj,j} \\ & = f'_j(\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j}) \psi^{m,J,j} \\ & \quad + (f'_j(\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j}) - f'_j(\mathbf{u}^{mj} \cdot G^{(m+1)J,j})) G^{mj,j} \\ & \quad + (f'_j(\mathbf{u}^{mj} \cdot G^{(m+1)J,j}) - f'_j(\mathbf{u}^{mj} \cdot G^{mj,j})) G^{mj,j} \\ & = f'_j(\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j}) \psi^{m,J,j} \\ & \quad + (d^{m,J} \cdot G^{(m+1)J,j}) G^{mj,j} \cdot \int_0^1 (1-t) \\ & \quad \times f''_j(\mathbf{u}^{mj} \cdot G^{(m+1)J,j} + t(d^{m,J} \cdot G^{(m+1)J,j})) dt \\ & \quad + (\mathbf{u}^{mj} \cdot \psi^{m,J,j}) G^{mj,j} \cdot \int_0^1 (1-t) \\ & \quad \times f''_j(\mathbf{u}^{mj} \cdot G^{mj,j} + t(\mathbf{u}^{mj} \cdot \psi^{m,J,j})) dt. \end{aligned} \quad (61)$$

Note (A2) and let  $\eta_c > 0$  be an upper bound of  $\{\eta_m\}_{m=0}^{\infty}$ . It follows from (36)–(38) that

$$\begin{aligned} & |\mathbf{u}^{mj} \cdot G^{(m+1)J,j} + t(d^{m,J} \cdot G^{(m+1)J,j})| \\ & \leq (\|\mathbf{u}^{mj}\| + \|d^{m,J}\|) \|G^{(m+1)J,j}\| \\ & \leq (C_2 + C_4 \eta_c) C_3, \quad t \in (0, 1), \end{aligned} \quad (62)$$

$$\begin{aligned} & |\mathbf{u}^{mj} \cdot G^{mj,j} + t(\mathbf{u}^{mj} \cdot \psi^{m,J,j})| \\ & \leq \|\mathbf{u}^{mj}\| (\|G^{mj,j}\| + \|\psi^{m,J,j}\|) \\ & \leq C_2(C_3 + C_5 \eta_c) = D_2 + C_2 C_5 \eta_c, \quad t \in (0, 1). \end{aligned} \quad (63)$$

According to (62), (63) and the proof of Lemma 4.1, there are positive constants  $C_{10}, C_{11} > 0$  such that

$$\left| \int_0^1 (1-t) f''_j(\mathbf{u}^{mj} \cdot G^{(m+1)J,j} + t d^{m,J} \cdot G^{(m+1)J,j}) dt \right| \leq C_{10}, \quad (64)$$

$$\left| \int_0^1 (1-t) f''_j(\mathbf{u}^{mj} \cdot G^{mj,j} + t \mathbf{u}^{mj} \cdot \psi^{m,J,j}) dt \right| \leq C_{11}. \quad (65)$$

By (43), we obtain  $|\mathbf{u}^{(m+1)J} \cdot G^{(m+1)J,j}| \leq C_2 C_3 = D_2$ . Employing (4) and (37), (38), (44), (64) and (65), and summing (61) from  $j = 1$  to  $j = J$ , we conclude that

$$\begin{aligned} & \|E_{\mathbf{u}}(\mathbf{w}^{(m+1)J})\| - \|E_{\mathbf{u}}(\mathbf{w}^{mj})\| \leq \|E_{\mathbf{u}}(\mathbf{w}^{(m+1)J}) - E_{\mathbf{u}}(\mathbf{w}^{mj})\| \\ & \leq C_{10} J \max_{\substack{1 \leq j \leq J \\ m \in \mathbb{N}}} (\|G^{(m+1)J,j}\| \|G^{mj,j}\|) \|d^{m,J}\| \\ & \quad + \left( J C_{4,1} + C_{11} J \max_{\substack{1 \leq j \leq J \\ m \in \mathbb{N}}} (\|\mathbf{u}^{mj}\| \|G^{mj,j}\|) \right) \max_{\substack{1 \leq j \leq J \\ m \in \mathbb{N}}} \|\psi^{m,J,j}\| \\ & \leq C_{12} \eta_m, \end{aligned} \quad (66)$$

where  $C_{12} = J C_3^2 C_4 C_{10} + J C_{4,1} C_5 + J D_2 C_5 C_{11}$ . Combining (59), (66) and Lemma 4.2 results in  $\lim_{m \rightarrow \infty} \|E_{\mathbf{u}}(\mathbf{w}^{mj})\| = 0$ .

Similarly as in the proof to (66), there exists a positive constant  $C_{13}$  such that

$$\|E_{\mathbf{u}}(\mathbf{w}^{mj+j}) - E_{\mathbf{u}}(\mathbf{w}^{mj})\| \leq C_{13} \eta_m. \quad (67)$$

Since

$$\begin{aligned} \|E_{\mathbf{u}}(\mathbf{w}^{mj+j})\| & \leq \|E_{\mathbf{u}}(\mathbf{w}^{mj+j}) - E_{\mathbf{u}}(\mathbf{w}^{mj})\| + \|E_{\mathbf{u}}(\mathbf{w}^{mj})\| \\ & \leq C_{13} \eta_m + \|E_{\mathbf{u}}(\mathbf{w}^{mj})\|, \end{aligned} \quad (68)$$

we have  $\lim_{m \rightarrow \infty} \|E_{\mathbf{u}}(\mathbf{w}^{mj+j})\| = 0$  for  $j = 1, 2, \dots, J$ . Similarly, we deduce that  $\lim_{m \rightarrow \infty} \|E_{\mathbf{v}_i}(\mathbf{w}^{mj+j})\| = 0$  for  $i = 1, \dots, n$ ,  $j = 1, 2, \dots, J$ , and

$$\lim_{m \rightarrow \infty} \|E_{\mathbf{w}}(\mathbf{w}^{mj+j})\| = 0, \quad j = 1, 2, \dots, J. \quad (69)$$

This immediately gives

$$\lim_{m \rightarrow \infty} \|E_{\mathbf{w}}(\mathbf{w}^m)\| = 0. \quad \square \quad (70)$$

**Proof of (21).** According to (A3), the sequence  $\{\mathbf{w}^m\}$  ( $m \in \mathbb{N}$ ) has a subsequence  $\{\mathbf{w}^{m_k}\}$  ( $k \in \mathbb{N}$ ) that is convergent to, say,  $\mathbf{w}^* \in \Omega_0$ . It follows from (20) and the continuity of  $E_{\mathbf{w}}(\mathbf{w})$  that

$$\|E_{\mathbf{w}}(\mathbf{w}^*)\| = \lim_{k \rightarrow \infty} \|E_{\mathbf{w}}(\mathbf{w}^{m_k})\| = \lim_{m \rightarrow \infty} \|E_{\mathbf{w}}(\mathbf{w}^m)\| = 0. \quad (71)$$

This implies that  $\mathbf{w}^*$  is a stationary point of  $E(\mathbf{w})$ . Hence,  $\{\mathbf{w}^m\}$  has at least one accumulation point and every accumulation point must be a stationary point.

Next, by reduction to absurdity, we prove that  $\{\mathbf{w}^m\}$  has precisely one accumulation point. Let us assume to the contrary that  $\{\mathbf{w}^m\}$  has at least two accumulation points  $\bar{\mathbf{w}} \neq \tilde{\mathbf{w}}$ . We write  $\mathbf{w}^m = (w_1^m, w_2^m, \dots, w_{n(p+1)}^m)^T$ . It is easy to see from (9)–(12) that  $\lim_{m \rightarrow \infty} \|\mathbf{w}^{m+1} - \mathbf{w}^m\| = 0$ , or equivalently,  $\lim_{m \rightarrow \infty} |w_i^{m+1} - w_i^m| = 0$  for  $i = 1, 2, \dots, n(p+1)$ . Without loss of generality, we assume that the first components of  $\bar{\mathbf{w}}$  and  $\tilde{\mathbf{w}}$  do not equal to each other, that is,  $\bar{w}_1 \neq \tilde{w}_1$ . For any real number  $\lambda \in (0, 1)$ , let  $w_1^\lambda = \lambda \bar{w}_1 + (1-\lambda)\tilde{w}_1$ . By Lemma 4.3, there exists a subsequence  $\{w_1^{m_{k_1}}\}$  of  $\{w_1^m\}$  converging to  $w_1^\lambda$  as  $k_1 \rightarrow \infty$ . Due to the boundedness of  $\{w_2^{m_{k_1}}\}$ , there is a convergent subsequence  $\{w_2^{m_{k_2}}\} \subset \{w_2^{m_{k_1}}\}$ . We define  $w_2^\lambda = \lim_{k_2 \rightarrow \infty} w_2^{m_{k_2}}$ . Repeating this procedure, we end up with decreasing subsequences  $\{m_{k_1}\} \supset \{m_{k_2}\} \supset \dots \supset \{m_{k_{n(p+1)}}\}$  with  $w_i^\lambda = \lim_{k_i \rightarrow \infty} w_i^{m_{k_i}}$  for each  $i = 1, 2, \dots, n(p+1)$ . Write  $\mathbf{w}^\lambda = (w_1^\lambda, w_2^\lambda, \dots, w_{n(p+1)}^\lambda)^T$ . Then, we see that  $\mathbf{w}^\lambda$  is an accumulation point of  $\{\mathbf{w}^m\}$  for any  $\lambda \in (0, 1)$ . But this means that  $\Omega_{0,1}$  has interior points, which contradicts (A4). Thus,  $\mathbf{w}^*$  must be a unique accumulation point of  $\{\mathbf{w}^m\}_{m=0}^\infty$ . This completes the proof of the strong convergence.  $\square$

#### 4.2. Convergence analysis for OGM-SS

Now, let the sequence  $\{\mathbf{w}^{mj+j}\}$  ( $m \in \mathbb{N}, j = 1, 2, \dots, J$ ) be generated by (14) and (15), and let

$$R^{m,j} = \Delta_j^m \mathbf{u}^{mj+j} - \Delta_j^m \mathbf{u}^{mj}, \quad (72)$$

$$r_i^{m,j} = \Delta_j^m \mathbf{v}_i^{mj+j} - \Delta_j^m \mathbf{v}_i^{mj}, \quad (73)$$

$$d^{m,l} = \mathbf{u}^{mj+l} - \mathbf{u}^{mj} = \sum_{j=1}^l \Delta_j^m \mathbf{u}^{mj+j} = \sum_{j=1}^l \Delta_j^m \mathbf{u}^{mj} + \sum_{j=1}^l R^{m,j}, \quad (74)$$

$$h_i^{m,l} = \mathbf{v}_i^{mj+l} - \mathbf{v}_i^{mj} = \sum_{j=1}^l \Delta_j^m \mathbf{v}_i^{mj+j} = \sum_{j=1}^l \Delta_j^m \mathbf{v}_i^{mj} + \sum_{j=1}^l r_i^{m,j}, \quad (75)$$

$$\psi^{m,l,j} = G^{mj+l, m, j} - G^{mj, m, j}, \quad (76)$$

$$m \in \mathbb{N}, j = 1, 2, \dots, J, l = 1, 2, \dots, J, i = 1, 2, \dots, n.$$

It is obvious that Lemmas 4.1–4.3 are not influenced by the new definitions. In place of Lemmas 4.4 and 4.5, we now have the following two Lemmas.

**Lemma 4.6.** *Let conditions (A1) and (A3) be valid, and let the sequence  $\{\mathbf{w}^{mj+j}\}$  be generated by (14) and (15). Then, there hold the following estimates with the same constants  $C_3$ – $C_7$  as in Lemma 4.4:*

$$\|G^{mj+j, m, k}\| \leq C_3, \quad (77)$$

$$\|d^{m,l}\| \leq C_4 \eta_m, \quad \|\psi^{m,l,j}\| \leq C_5 \eta_m, \quad (78)$$

$$\|R^{m,j}\| \leq C_6 \eta_m^2, \quad \|r_i^{m,j}\| \leq C_7 \eta_m^2, \quad (79)$$

where  $m \in \mathbb{N}; j, k = 1, 2, \dots, J; l = 1, 2, \dots, J; i = 1, 2, \dots, n$ .

**Proof.** According to (36), we have

$$|\mathbf{v}_i^{mj+j} \cdot \mathbf{x}^{m,k}| \leq \|\mathbf{v}_i^{mj+j}\| \max_{1 \leq k \leq J} \|\mathbf{x}^k\| \leq C_1 C_2 \equiv D_1. \quad (80)$$

Thus, there exists a positive constant  $C_{3,1}$  such that

$$\max_{|t| \leq D_1} |g(t)| = C_{3,1}, \quad (81)$$

$$\|G^{mj+j, m, k}\| = \|G(\mathbf{v}^{mj+j} \mathbf{x}^{m,k})\| \leq \sqrt{n} C_{3,1} \equiv C_3. \quad (82)$$

Similarly, (78) and (79) can be proved after adjusting the corresponding superscripts in the proof to Lemma 4.4.  $\square$

**Lemma 4.7.** *Let the sequence  $\{\mathbf{w}^{mj+j}\}$  be generated by (14) and (15). Under assumptions (A1) and (A3), there holds*

$$E(\mathbf{w}^{(m+1)J}) \leq E(\mathbf{w}^{mj}) - \eta_m \|E_{\mathbf{w}}(\mathbf{w}^{mj})\|^2 + C_8 \eta_m^2, \quad (83)$$

$$(m = 0, 1, \dots)$$

where  $C_8 > 0$  is the same constant defined in Lemma 4.5.

**Proof.** As in the proof to Lemma 4.6, we only need to adjust some superscripts. For example, corresponding to (54), we change the related superscripts and get

$$f_j'(\mathbf{u}^{mj} \cdot G^{mj, m, j}) \mathbf{u}^{mj} \cdot \psi^{m, J, j}$$

$$= f_j'(\mathbf{u}^{mj} \cdot G^{mj, m, j}) \sum_{i=1}^n u_i^{mj} g'(v_i^{mj} \cdot \mathbf{x}^{m, j}) h_i^{m, J} \cdot \mathbf{x}^{m, j}$$

$$+ f_j'(\mathbf{u}^{mj} \cdot G^{mj, m, j}) \sum_{i=1}^n u_i^{mj} (h_i^{m, J} \cdot \mathbf{x}^{m, j})^2$$

$$\times \int_0^1 (1-t) g''(v_i^{mj} \cdot \mathbf{x}^{m, j} + t(h_i^{m, J} \cdot \mathbf{x}^{m, j})) dt. \quad (84)$$

The details are left to the interested readers.  $\square$

**Proof of Theorem 3.1 for OGM-SS.** We can use Lemmas 4.1–4.3 and Lemmas 4.6–4.7 to obtain the weak and strong convergence results for OGM-SS precisely as in the proof to Theorem 3.1 for OGM-F.  $\square$

## 5. Conclusions

In this paper, we present a comprehensive study on the weak and strong convergence for three-layer BP neural networks. Compared with the existing convergence results, the corresponding assumptions are more relaxed. Our convergence analysis holds for more extensive BP neural networks, e.g., S–S, S–P, P–P and P–S type neural networks.

## References

- Bertsekas, D. P., & Tsitsiklis, J. N. (1996). *Neuro-dynamic programming*. Athena Scientific.
- Chakraborty, D., & Pal, N. R. (2003). A novel training scheme for multilayered perceptrons to realize proper generalization and incremental learning. *IEEE Transactions on Neural Networks*, 14, 1–14.
- Fine, T. L., & Mukherjee, S. (1999). Parameter convergence and learning curves for neural networks. *Neural Computation*, 11, 747–769.
- Finnoff, W. (1994). Diffusion approximations for the constant learning rate backpropagation algorithm and resistance to local minima. *Neural Computation*, 6, 285–295.
- Heskes, T., & Wiegerinck, W. (1996). A theoretical comparison of batch-mode, on-line, cyclic, and almost-cyclic learning. *IEEE Transactions on Neural Networks*, 7, 919–925.
- LeCun, Y. (1985). Une procedure d'apprentissage pour reseau à seuil asymmetrique. *A la Frontieredel'Intelligence Artificielle des Sciences de la Connaissance des Neurosciences*, 85, 599–604.
- Li, Z. X., & Ding, X. S. (2005). Prediction of stock market by bp neural networks with technical indexes as input. *Numerical Mathematics: A Journal of Chinese Universities*, 27, 373–377.
- Li, Z. X., Wu, W., & Tian, Y. L. (2004). Convergence of an online gradient method for feedforward neural networks with stochastic inputs. *Journal of Computational and Applied Mathematics*, 163, 165–176.
- Liang, Y. C., Feng, D. P., Lee, H. P., Lim, S. P., & Lee, K. H. (2002). Successive approximation training algorithm for feedforward neural networks. *Neurocomputing*, 42, 11–322.
- Mangasarian, O. L., & Solodov, M. V. (1994). Serial and parallel backpropagation convergence via nonmonotone perturbed minimization. *Optimization Methods and Software*, 4, 117–134.
- Nakama, T. (2009). Theoretical analysis of batch and on-line training for gradient descent learning in neural networks. *Neurocomputing*, 73, 151–159.
- Parker, D. B. (1982). Learning-logic, invention report. Stanford University, Stanford, Calif.



- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagation errors. *Nature*, 323, 533–536.
- Shao, H. M., Wu, W., & Liu, W. B. (2007). Convergence of BP algorithm for training MLP with linear output. *Numerical Mathematics: A Journal of Chinese Universities (English Series)*, 16, 193–202.
- Tadic, V., & Stankovic, S. (2000). Learning in neural networks by normalized stochastic gradient algorithm: local convergence. In *Proceedings of the 5th seminar neural networks application electronic engineering*.
- Terence, D. S. (1989). Optimal unsupervised learning in a single-layer linear feedforward neural network. *Neural Networks*, 2, 459–473.
- Werbos, P. J. (1974). Beyond regression: new tools for prediction and analysis in the behavioral sciences. Ph.D. thesis. Harvard University, Cambridge, MA.
- Wilson, D. R., & Martinez, T. R. (2003). The general inefficiency of batch training for gradient descent learning. *Neural Networks*, 16, 1429–1451.
- Wu, W., & Xu, Y. S. (2002). Deterministic convergence of an on-line gradient method for neural networks. *Journal of Computational and Applied Mathematics*, 144, 335–347.
- Wu, W., Feng, G. R., Li, Z. X., & Xu, Y. S. (2005). Deterministic convergence of an online gradient method for BP neural networks. *IEEE Transactions on Neural Networks*, 16, 533–540.
- Wu, W., Feng, G. R., & Li, X. (2002). Training multilayer perceptrons via minimization of sum of ridge functions. *Advances in Computational Mathematics*, 17, 331–347.
- Wu, W., & Shao, Z. Q. (2003). Convergence of online gradient methods for continuous perceptrons with linearly separable training patterns. *Applied Mathematics Letters*, 16, 999–1002.
- Wu, W., Shao, H. M., & Qu, D. (2005). Strong convergence of gradient methods for BP networks training. In *Proceedings of 2005 international conference on neural networks and brains* (pp. 332–334).
- Xu, Z. B., Zhang, R., & Jin, W. F. (2009). When on-line BP training converges. *IEEE Transactions on Neural Networks*, 20, 1529–1539.
- Zhang, H. S., Wu, W., Liu, F., & Yao, M. C. (2009). Boundedness and convergence of online gradient method with penalty for feedforward neural networks. *IEEE Transactions on Neural Networks*, 20, 1050–1054.