
Convergence of Online Gradient Method for Pi-sigma Neural Networks with Inner-penalty Terms

Kh. Sh. Mohamed^{1,2}, Xiong Yan³, Y. Sh. Mohammed^{4,5}, Abd-Elmoniem A. Elzain^{5,6}, Habtamu Z. A.², Abdrhaman M. Adam²

¹Mathematical Department, College of Science, Dalanj University, Dalanj, Sudan

²School of Mathematical Sciences, Dalian University of Technology, Dalian, China

³School of Science, Liaoning University of Science & Technology, Anshan, China

⁴Physics Department, College of Education, Dalanj University, Dalanj, Sudan

⁵Department of Physics, College of Science & Art, Qassim University, Oklat Al- Skoor, Saudi Arabia

⁶Department of Physics, University of Kassala, Kassala, Sudan

Email address:

khshm7@yahoo.com (Kh. Sh. Mohamed), xy-zhxw@163.com (Xiong Yan), yshm@yahoo.com (Y. Sh. Mohammed), Abdelmoniem1@yahoo.com (Abd-Elmoniem A. E.), habtamuz@mail.dlut.edu.cn (Habtamu Z. A.), abdelrhaman013@yahoo.com (A. M. Adam)

To cite this article:

Kh. Sh. Mohamed, Xiong Yan, Y. Sh. Mohammed, Abd-Elmoniem A. Elzain, Habtamu Z. A., Abdrhaman M. Adam. Convergence of Online Gradient Method for Pi-sigma Neural Networks with Inner-penalty Terms. *American Journal of Neural Networks and Applications*. Vol. 2, No. 1, 2016, pp. 1-5. doi: 10.11648/j.ajjna.20160201.11

Received: March 14, 2016; **Accepted:** March 30, 2016; **Published:** May 10, 2016

Abstract: This paper investigates an online gradient method with inner-penalty for a novel feed forward network it is called pi-sigma network. This network utilizes product cells as the output units to indirectly incorporate the capabilities of higher-order networks while using a fewer number of weights and processing units. Penalty term methods have been widely used to improve the generalization performance of feed forward neural networks and to control the magnitude of the network weights. The monotonicity of the error function and weight boundedness with inner-penalty term and both weak and strong convergence theorems in the training iteration are proved.

Keyword: Convergence, Pi-sigma Network, Online Gradient Method, Inner-penalty, Boundedness

1. Introduction

A novel higher order feedforward polynomial neural network is known to provide inherently more powerful mapping abilities than traditional feed forward neural network called the pi-sigma network (PSN) [2]. This network utilizes product cells as the output units to indirectly incorporate the capabilities of higher-order networks while using a fewer number of weights and processing units. The neural networks consisting of the PSN modules has been used effectively in pattern classification and approximation problems [1, 4, 10, 11]. There are two ways of training to updating weight: The first approach, batch (offline) training [18], the weights are modified after each training pattern is presented to the network. Second approach, online

training, the weights updating immediately after each training sample is fed see [13]. The penalty term is often introduced into the network training algorithms has been widely used so as to control the magnitude of the weights and to improve the generalization performance of the network [6, 8], here the generalization performance refers to the capacity of a neural network to give correct outputs for untrained data. Specially cause, in the second approach the training weights updating become very large and over-fitting tends to occur, by adding the penalty term in into the cost function, when use second approach has been successfully application see [3, 7, 12, 14], which acts as a brute-force to drive unnecessary weights to zero and to prevent the weights from taking too large in the training process. In the work area of penalty term at the same of the inner-penalty term (IP), which have worked to reduce the magnitude of the network weights with efficiency

improve the generalization performance of the network [5, 9, 17]. In this paper, we prove the (strong and weak) convergence of the online gradient with inner penalty and the monotonicity of the error function and the weight sequence are uniformly bounded during the training procedure with inner-penalty.

The rest of this paper is organized as follows. The neural network structure and the online gradient method with inner-penalty are described in Section 2. The preliminary lemmas are disruption in Section 3. The convergence results are presented and the rigorous proofs of the main results are provided in Section 4. Finally, in Section 5 we conclusions this study.

2. PSN-IPAlgorithm

PSN is a higher order feed forward polynomial neural network consisting of a single hidden layer. The hidden layer has summing units where as the output layer has product units. PSN, which has a three-layer network consisting of p input units, N summation units, and 1 product layers. Let $\omega_k = (\omega_{k1}, \omega_{k2}, \dots, \omega_{kp})^T \in \mathbb{R}^p$ ($1 \leq k \leq N$) the weight vectors connecting the input and summing units, and write $\omega = (\omega_1^T, \omega_2^T, \dots, \omega_N^T) \in \mathbb{R}^{Np}$. We have included a special input unit ξ_p , corresponding to the biases ω_{kp} , with fixed value-1. The structure of PSN is shown in Figure 1.

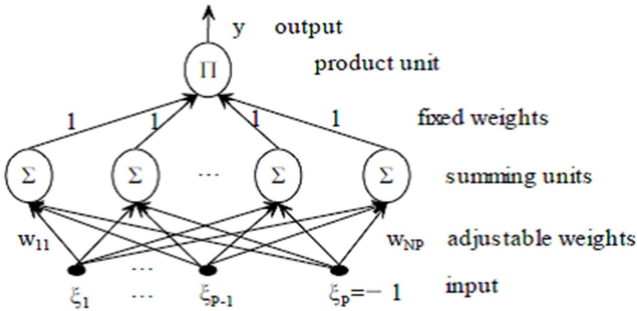


Figure 1. PSN structure with a single output

Where $\xi^j = (\xi_1^j, \dots, \xi_{p-1}^j, \xi_p^j) \in \mathbb{R}^p$ and $\xi_p^j \equiv -1$. Assume $g: \mathbb{R} \rightarrow \mathbb{R}$ is a given activation function. For an input $\xi^j \in \mathbb{R}^p$, the output of the network is

$$y = g(\prod_{i=1}^N (\omega_i \cdot \xi^j)) \quad (1)$$

The network supplied with a given set of training samples $\{\xi^j, o^j\}_{j=1}^J \subset \mathbb{R}^p \times \mathbb{R}$. The error function with a inner penalty given by

$$\begin{aligned} E(\omega) &= \frac{1}{2} \sum_{j=1}^J (o^j - g(\prod_{i=1}^N (\omega_i \cdot \xi^j)))^2 + \frac{\lambda}{2} \sum_{j=1}^J \sum_{k=1}^N (\omega_k \cdot \xi^j)^2 \\ &= \sum_{j=1}^J g_j(\prod_{i=1}^N (\omega_i \cdot \xi^j)) + \frac{\lambda}{2} \sum_{j=1}^J \sum_{k=1}^N (\omega_k \cdot \xi^j)^2 \end{aligned} \quad (2)$$

Where $\lambda > 0$ is a inner penalty coefficient and $g_j(t) =$

$\frac{1}{2} (o^j - g(t))^2$ The gradient function is given by

$$E_{\omega_k}(\omega) = \sum_{j=1}^J [g'_j(\prod_{i=1}^N (\omega_i \cdot \xi^j)) \prod_{i=1, i \neq k}^N (\omega_i \cdot \xi^j) + \lambda (\omega_k \cdot \xi^j)] \xi^j \quad (3)$$

Given an initial weight ω^0 , the online method with inner penalty updates them iteratively by the form

$$\omega_k^{m+1} = \omega_k^m + \Delta \omega_k^m, \quad m = 0, 1, \dots \quad (4)$$

$$\Delta \omega_k^m = [g'_j(\prod_{i=1}^N (\omega_i^m \cdot \xi^j)) \prod_{i=1, i \neq k}^N (\omega_i^m \cdot \xi^j) + \lambda (\omega_k^m \cdot \xi^j)] \xi^j \quad (5)$$

Where $\eta > 0$ is the learning rate in the m th training cycle. We denote by $\|\cdot\|$ the usual Euclidean norm and the corresponding derived matrix norm and the following Assumption is imposed throughout this paper.

Assumption 1.

$$|g_j(t)|, |g'_j(t)|, |g''_j(t)| \leq M_1, \forall t \in \mathbb{R}, 1, 1 \leq j \leq J$$

Assumption 2.

$$\|\xi^j\| \& |\omega_k^m \cdot \xi^j| < M_1, 1 \leq k \leq N, 1 \leq j \leq J, m = 0, 1, \dots$$

Assumption 3.

The learning rate η and penalty parameter λ are chosen to satisfy the condition: $0 \leq \eta \leq 1/(\lambda \tilde{M} + M)$

Assumption 4.

$\{\omega_k^m\}_{m=0,1,\dots}$ are contained in bounded closed region $\Theta \subset \mathbb{R}^{Np}$, and there are exist points in set $\Theta_0 = \{\omega \in \Theta | E_{\omega}(\omega) = 0\}$.

3. Preliminary Lemmas

The next lemma present the monotonicity of the sequence $\{E(\omega)\}$. It is essential for the proof of weakly convergence of PSN with penalty, presented in the following Theorems. For sake of description, we denote

$$r_{k,j}^m = \Delta \omega_k^{m+1} - \Delta \omega_k^m \quad (6)$$

$$\psi_j^{m+1} = \prod_{i=1}^N (\omega_i^{m+1} \cdot \xi^j) \quad (7)$$

$$\varphi_{k,j}^{m+1} = \prod_{i=1, i \neq k}^N (\omega_i^{m+1} \cdot \xi^j) \quad (8)$$

$$1 \leq j \leq J, 1 \leq i \leq N, m = 0, 1, \dots$$

To begin with, first we present a few lemmas as preparation to prove Theorems

Lemma 1. Let Assumption 1-2 are valid, there hold

$$(i) \|\psi_j^{m+1} - \psi_j^m\| \leq M_2 \sum_{k=1}^N \|\Delta \omega_k^m\| \quad (9)$$

$$(ii) \|\varphi_{k,j}^{m+1} - \varphi_{k,j}^m\| \leq \tilde{M}_2 \sum_{k=1}^N \|\Delta \omega_k^m\| \quad (10)$$

Proof. By Assumption 2 and Cauchy- Schwartz inequality, we have

$$|\psi_j^{m+1} - \psi_j^m| \leq \left| \prod_{i=1}^{N-1} (\omega_i^{m+1} \cdot \xi^j) \right| |(\omega_N^{m+1} - \omega_N^m) \xi^j|$$

$$\begin{aligned}
 & + \left| \prod_{i=1}^{N-2} (\omega_i^{m+1} \cdot \xi^j) (\omega_N^m \cdot \xi^j) \right| (\omega_{N-1}^{m+1} - \omega_{N-1}^m) \xi^j \\
 & + \dots + \left| \prod_{i=2}^N (\omega_i^m \cdot \xi^j) \right| |(\omega_1^{m+1} - \omega_1^m) \xi^j| \\
 & \leq M_1^{N-1} \|\xi^j\| \sum_{k=1}^N \|\Delta \omega_k^m\| \leq M_1^N \sum_{k=1}^N \|\Delta \omega_k^m\| \\
 & \leq M_2 \sum_{k=1}^N \|\Delta \omega_k^m\| \tag{11}
 \end{aligned}$$

where $M_2 = M_1^N$, $1 \leq j \leq J$, $1 \leq k \leq N$, $m = 0, 1, \dots$. Similarly, we get

$$\|\varphi_{k,j}^{m+1} - \varphi_{k,j}^m\| \leq \tilde{M}_2 \sum_{k=1}^N \|\Delta \omega_k^m\| \tag{12}$$

Here, $\tilde{M}_2 = M_1^{N-1}$. This completes the proof.

Lemma 2. Suppose Assumptions 1~3 are satisfied, and the weight sequence $\{\omega_k^m\}_{m=0,1,\dots}$ is generated by (4) ~ (5), then

$$\begin{aligned}
 & E(\omega_k^{m+1}) - E(\omega_k^m) \\
 & \leq -\left(\frac{1}{\eta} - \lambda - M_3 - M_4 - M_5\right) \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{13}
 \end{aligned}$$

Proof. Applying Taylor's formula to extend $g_j(\psi_j^{m+1})$ at (ψ_j^m) , we have

$$\begin{aligned}
 & g_j(\psi_j^{m+1}) - g_j(\psi_j^m) = g_j'(\psi_j^m) (\varphi_{k,j}^m) \xi^j \sum_{k=1}^N (\omega_k^{m+1} - \omega_k^m) \\
 & + \frac{1}{2} g_j''(t_1) (\psi_j^{m+1} - \psi_j^m)^2 \\
 & + \frac{1}{2} \sum_{\substack{k_1, k_2=1 \\ k_1 \neq k_2}}^N \left(\prod_{\substack{i=1 \\ i \neq k_1, k_2}}^N t_2 \right) (\omega_{k_1}^{m+1} - \omega_{k_1}^m) (\omega_{k_2}^{m+1} - \omega_{k_2}^m) (\xi^j)^2 \tag{14}
 \end{aligned}$$

where $t_1, t_2 \in \mathbb{R}$ are on the line segment between ψ_j^{m+1} and ψ_j^m . After dealing with (14) by accumulation $g_j(\psi_j^{m+1})$ for $1 \leq j \leq J$, we obtain from (2), (4), (5) and Taylor's formula we obtain

$$\begin{aligned}
 & E(\omega_k^{m+1}) - E(\omega_k^m) = \sum_{j=1}^J [g_j(\psi_j^{m+1}) - g_j(\psi_j^m)] \\
 & + \frac{\lambda}{2} \sum_{j=1}^J \sum_{k=1}^N [(\omega_k^{m+1} \cdot \xi^j)^2 - (\omega_k^m \cdot \xi^j)^2] \\
 & = \sum_{j=1}^J [g_j'(\psi_j^m) (\varphi_{k,j}^m) \xi^j + \lambda (\omega_k^m \cdot \xi^j) \xi^j] (\Delta \omega_k^m) \\
 & + \lambda \sum_{j=1}^J \sum_{k=1}^N (\Delta \omega_k^m \cdot \xi^j)^2 + \mathfrak{S}_1 + \mathfrak{S}_2 + \mathfrak{S}_3 + \mathfrak{S}_4 \\
 & = -\frac{1}{\eta} \sum_{k=1}^N \|\Delta \omega_k^m\|^2 + \mathfrak{S}_1 + \mathfrak{S}_2 + \mathfrak{S}_3 + \mathfrak{S}_4 \tag{15}
 \end{aligned}$$

where

$$|\mathfrak{S}_1| \leq \frac{1}{\eta} \sum_{j=1}^J (\Delta \omega_k^m) \cdot \sum_{k=1}^N (r_{k,j}^m) \tag{16}$$

$$|\mathfrak{S}_2| \leq \frac{1}{2} \sum_{j=1}^J g_j''(t_1) (\psi_j^{m+1} - \psi_j^m)^2 \tag{17}$$

$$|\mathfrak{S}_3| \leq \frac{1}{2} \sum_{j=1}^J g_j'(\psi_j^m) \sum_{\substack{k_1, k_2=1 \\ k_1 \neq k_2}}^N \left(\prod_{\substack{i=1 \\ i \neq k_1, k_2}}^N t_2 \right) (\Delta \omega_{k_1}^m \cdot \Delta \omega_{k_2}^m) (\xi^j)^2 \tag{18}$$

$$|\mathfrak{S}_4| \leq \frac{\lambda}{2} \sum_{j=1}^J \sum_{k=1}^N (\Delta \omega_k^m \cdot \xi^j)^2 \tag{19}$$

By Assumption 1, (4) ~ (5), Lemma 1 and the mean value theorem gives

$$\begin{aligned}
 \|r_{k,j}^m\| & = \|g_j'(\psi_j^{m+1}) \varphi_{k,j}^{m+1} - g_j'(\psi_j^m) \varphi_{k,j}^m\| \|\xi^j\| \\
 & + \lambda \|\omega_k^{m+1} - \omega_k^m\| \|\xi^j\|^2 \\
 & \leq \|g_j''(t_1) (\psi_j^{m+1} - \psi_j^m) (\varphi_{k,j}^{m+1}) \xi^j\| \\
 & + \|g_j'(\psi_j^m) (\varphi_{k,j}^{m+1} - \varphi_{k,j}^m)\| + \lambda \|\omega_k^{m+1} - \omega_k^m\| \|\xi^j\|^2 \\
 & \leq (M_2 M_1^{N+1} + \tilde{M}_2 M_1^2 + \lambda M_1^2) \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{20}
 \end{aligned}$$

Thus with (16) gives

$$\begin{aligned}
 |\mathfrak{S}_1| & \leq \frac{1}{\eta} \sum_{k=1}^N \|\Delta \omega_k^m \cdot r_{k,j}^m\| \\
 & \leq M_3 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{21}
 \end{aligned}$$

Here, $M_3 = M_2 M_1^{N+1} + \tilde{M}_2 M_1^2 + \lambda M_1^2 + 1$. By (4), (9) in Lemma 1 and Cauchy- Schwartz inequality, we have

$$\begin{aligned}
 |\mathfrak{S}_2| & \leq \frac{1}{2} M_1 |\psi_j^{m+1} - \psi_j^m|^2 \\
 & \leq \left(\frac{1}{2} M_2^2 + \lambda M_3\right) M_1 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \\
 & \leq M_4 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{22}
 \end{aligned}$$

where $M_4 = \frac{1}{2} M_1 M_2^2 + \lambda M_1^2$. It follows from Assumption 1~2, (2) and Taylor's formula, we obtain

$$\begin{aligned}
 |\mathfrak{S}_3| & \leq \frac{1}{2} M_1^{N+1} \sum_{j=1}^J \sum_{\substack{k_1, k_2=1 \\ k_1 \neq k_2}}^N \|\Delta \omega_{k_1}^m \cdot \Delta \omega_{k_2}^m\| \\
 & \leq \frac{1}{2} M_1^{N+1} J(N-1) \sum_{j=1}^J \sum_{k=1}^N \|\Delta \omega_{k_1}^m \cdot \Delta \omega_{k_2}^m\| \\
 & \leq M_5 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{23}
 \end{aligned}$$

where $M_5 = \frac{1}{2} M_1^{N+1} J(N-1)$. By Assumption (2) and (4) leads

$$\begin{aligned}
 |\mathfrak{S}_4| & \leq \lambda \|\xi^j\|^2 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \leq \frac{1}{2} \lambda M_1^2 \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \\
 & \leq \lambda \tilde{M} \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \tag{24}
 \end{aligned}$$

where $\tilde{M} = \frac{1}{2}M_1^2$. Collate (20) ~ (24) into (15) gives

$$E(\omega_k^{m+1}) - E(\omega_k^m) \leq -\left(\frac{1}{\eta} - \lambda - M_3 - M_4 - M_5\right) \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \quad (25)$$

This completes the proof.

Lemma 3. Suppose that $F: \mathbb{R}^K \rightarrow \mathbb{R}$ is continuous and differentiable on a compact set $\tilde{D} \subset \mathbb{R}^K$ and that $\Theta = \{Z \in \tilde{D} | \nabla h(Z) = 0\}$ has only finite number of point. If a sequence $\{Z^m\}_{m=1}^\infty \in \tilde{D}$ satisfies then $\lim_{m \rightarrow \infty} \|Z^{m+1} - Z^m\| = 0$, $\lim_{m \rightarrow \infty} \|\nabla h(Z^m)\| = 0$. Then there exists a point $Z^* \in \Theta$ such that $\lim_{m \rightarrow \infty} Z^m = Z^*$.

Proof. This result is basically the same as Theorem 14.1.5 in [16], and the detailed proof is thus omitted.

4. Convergence Theorems

Now, we can elucidate and proofs the convergence theorems, which we needed

Theorem 1. (Monotonicity): Let Assumption 1~3 are valid and the weight sequence $\{\omega_k^m\}_{m=0,1,\dots}$ be generated by (4) ~ (5), then

$$E(\omega_k^{m+1}) - E(\omega_k^m), m = 0, 1, \dots \quad (26)$$

Proof. Let

$$M = M_3 + M_4 + M_5 \quad (27)$$

By Assumption 3, which satisfies

$$0 < \eta < \frac{1}{\lambda \tilde{M} + M} \quad (28)$$

Thus with (28) and Lemma 2, we have

$$E(\omega_k^{m+1}) - E(\omega_k^m) \leq -\left(\frac{1}{\eta} - \lambda \tilde{M} - M\right) \sum_{k=1}^N \|\Delta \omega_k^m\|^2 \leq 0 \quad (29)$$

This completes the proof of the Theorem 1.

Theorem 2. (Boundedness): Suppose that Assumption of Theorem 1 are valid, the weight sequence $\{\omega_k^m\}_{m=0,1,\dots}$ be generated by (4) ~ (5) are uniformly bounded.

Proof. By Assumption 1, and Theorem 1, we have

$$E(\omega_k^m) - E(\omega_k^{m-1}) \leq \dots \leq E(\omega_k^0) = \sum_{j=1}^J g_j(\psi_j^0) + \frac{\lambda}{2} \sum_{j=1}^J \sum_{k=1}^N (\omega_k^0 \cdot \xi^j)^2 \leq C \quad (30)$$

and

$$C = JM_1 \left(1 + \frac{\lambda}{2} \sum_{k=1}^N \|\omega_k^0\|^2\right) \quad (31)$$

From (2), (30) gives

$$\lambda (\omega_k^m \cdot \xi^j)^2 \leq 2E(\omega_k^m) \leq 2C, j = 1, 2, \dots, J \quad (32)$$

By (4) ~ (5), we have

$$\omega_k^m = \omega_k^0 - \eta \sum_{t=1}^{m-1} \sum_{k=1}^N [g'_j(\psi_j^t) \varphi_{k,j}^t + \lambda (\omega_k^t \cdot \xi^j)] \xi^j \quad (33)$$

Let the second part of above equation be ω_{k1}^m , Denote $\mathbb{R}_1 = \text{span}\{\xi^1, \xi^2, \dots, \xi^J\} \subset \mathbb{R}^n$ and $\mathbb{R}_2 = \mathbb{R}_1^\perp$ be the orthogonal complement space of \mathbb{R}_1 . Denote the second part of (33) by ω_{k2}^m , obviously $\omega_{k1}^m \in \mathbb{R}_1$. we divide ω_k^0 into $\omega_k^0 = \omega_{k1}^0 + \omega_{k2}^0$, where $\omega_{k1}^0 \in \mathbb{R}_1$ and $\omega_{k2}^0 \in \mathbb{R}_2$. Then $\omega_k^m = (\omega_{k1}^0 + \omega_{k1}^m) \oplus \omega_{k2}^0 = \tilde{\omega}_{k1}^m \oplus \omega_{k2}^0$. Applying this to (33) we have

$$|d_t| := |\tilde{\omega}_{k1}^m \cdot \xi^t| = |\omega_k^m \cdot \xi^j| \leq \sqrt{\frac{2C}{\lambda}}, t = 1, 2, \dots, T \quad (34)$$

Suppose $\{\xi^{i_1}, \xi^{i_2}, \dots, \xi^{i_T}\} (i_t \in \{1, \dots, J\}, t = 1, 2, \dots, T)$ is a base of the space \mathbb{R}_1 . There are $a_t \in \mathbb{R} (t = 1, 2, \dots, T)$ such that $\tilde{\omega}_{k1}^m = a_1 \xi^{i_1} + \dots + a_T \xi^{i_T}$. Then $(a_1 \xi^{i_1} + \dots + a_T \xi^{i_T}) \cdot \xi^{i_t} = d_t, t = 1, \dots, T$. we get

$$\begin{pmatrix} \xi^{j_1} \cdot \xi^{j_1} & \dots & \xi^{j_t} \cdot \xi^{j_1} \\ \vdots & \vdots & \vdots \\ \xi^{j_1} \cdot \xi^{j_t} & \dots & \xi^{j_t} \cdot \xi^{j_t} \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_T \end{pmatrix} = \begin{pmatrix} d_1 \\ \vdots \\ d_T \end{pmatrix} \quad (35)$$

Is a base, the coefficient determinant equal to zero, and the system of the linear equations has a unique solution. Assume that the coefficient determinant equals to. Then the solution is as follows

$$L = \begin{vmatrix} \xi^{j_1} \cdot \xi^{j_1} & \dots & \xi^{j_{t-1}} \cdot \xi^{j_1} & \dots & \xi^{j_1} \cdot \xi^{j_1} & \dots & \xi^{j_1} \cdot \xi^{j_1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \xi^{j_1} \cdot \xi^{j_t} & \dots & \xi^{j_{t-1}} \cdot \xi^{j_t} & \dots & \xi^{j_t} \cdot \xi^{j_1} & \dots & \xi^{j_t} \cdot \xi^{j_t} \end{vmatrix}$$

Then the solution is as follows

$$a_t = L \cdot D^{-1} \quad (36)$$

Let the maximum absolute value of all the sub-determinant with rank $(T - 1)$ of the coefficient determinant is D' , then $|a_t| \leq |D'| \cdot |D^{-1}| \cdot \sum_{t=0}^T |d_t|$. By (34) we have $|a_t| \leq |D'| \cdot |D^{-1}| \cdot T \cdot \sqrt{\frac{2C}{\lambda}}, t = 1, 2, \dots, T$. Denote $\tilde{C} = \max_{1 \leq t \leq T} \|\xi^{j_1}\|$, then

$$\begin{aligned} \|\tilde{\omega}_{k1}^m\| &= \|a_1 \xi^{j_1} + \dots + a_T \xi^{j_T}\| \\ &\leq |D'| \cdot |D^{-1}| \cdot \tilde{C} \cdot T^2 \cdot \sqrt{\frac{2C}{\lambda}} \end{aligned} \quad (37)$$

That is $\tilde{\omega}_{k1}^m$ are bounded uniformly bounded. So from (29), we know ω_k^m are uniformly bounded. In all, we get $\{\omega_k^m\}_{m=0,1,\dots}$ are uniformly bounded, i.e., there exist a bounded closed region $D \subset \mathbb{R}^k$ such that $\{\omega_k^m\} \subset D$.

Theorem 3. (Weak convergence): Suppose that Assumption 1~3 are valid and the weight sequence $\{\omega_k^m\}_{m=0,1,\dots}$ be generated by (4) ~ (5), then

$$\lim_{m \rightarrow \infty} \|E_{\omega_k}(\omega_k^m)\| = 0, \quad (38)$$

Furthermore, if Assumption 4 is also valid, we have the strong convergence: There exists $\omega^* \in \Theta_0$ such that

$$\lim_{m \rightarrow \infty} \omega_k^m = \omega_k^* \quad (39)$$

Proof. By (28) and setting $\beta > 0$ such that $\beta = \frac{1}{\eta} - \lambda\tilde{M} - M$, we have

$$E(\omega_k^{m+1}) \leq E(\omega_k^m) - \beta \sum_{k=1}^N \|\Delta\omega_k^m\|^2 \leq \dots \leq E(\omega_k^0) - \beta \sum_{i=0}^m \sum_{k=1}^N \|\Delta\omega_k^i\|^2 \quad (40)$$

Since $E(\omega_k^{m+1}) > 0$ for any $m = 0, 1, \dots$. We set $m \rightarrow \infty$

$$\sum_{i=0}^m \sum_{k=1}^N \|\Delta\omega_k^i\|^2 \leq \frac{E(\omega_k^0)}{\beta} < \infty \quad (41)$$

Combining (3) ~ (5), immediately gives the weak convergence result:

$$\lim_{m \rightarrow \infty} \|E_{\omega_k}(\omega_k^m)\| = 0, m = 0, 1, \dots \quad (42)$$

Next we prove the strong convergence it follows from (4)~(5) and (42) that leads

$$\lim_{m \rightarrow \infty} \|\Delta\omega_k^m\| = 0, \quad 0 \leq k \leq N \quad (43)$$

Note that the error function $E(\omega^m)$ defined in (2) is continuously differentiable. By (43), Assumptions 4 and Lemma 3, immediately get the desired result. This completes the proof.

5. Conclusion

Through our study of this paper, the monotonicity of the error function $E(\omega^m)$ in formula (2) and the weight sequence boundedness $\{\omega_k^m\}_{m=0,1,\dots}$ via formula (4) ~ (5) for the online gradient method with inner-penalty are presented, under those condition both weakly and strongly convergence theorems are proved.

Acknowledgment

We gratefully acknowledge Dalanj University for supporting this research. And our thanks to the anonymous reviewers and the editors for their previous helpful comments and valuable suggestion for their kind helps during the period of the research.

References

- [1] A J Hussaina and P Liatsisb, Recurrent pi-sigma networks for DPCM image coding. *Neurocomputing*, 55(2002) 363-382.
- [2] Y Shin and J Ghosh, The pi-sigma network: An efficient higher-order neural network for pattern classification and function approximation. *International Joint Conference on Neural Networks*, 1(1991) 13-18.
- [3] P L Bartlett, For valid generalization, the size of the weights is more important than the size of the network, *Advances in Neural Information Processing Systems* 9 (1997) 134-140.
- [4] L J Jiang, F Xu and S R Piao, Application of pi-sigma neural network to real-time classification of seafloor sediments. *Applied Acoustics*, 24(2005) 346-350.
- [5] R Reed, Pruning algorithms-a survey. *IEEE Transactions on Neural Networks* 8 (1997) 185-204.
- [6] G Hinton, Connectionist learning procedures, *Artificial Intelligence* 40(1989)185-243.
- [7] S Geman, E Bienenstock, R Doursat, Neural networks and the bias/variance dilemma, *Neural Computation* 4 (1992) 1-58.
- [8] S Loone and G Irwin, Improving neural network training solutions using regularisation, *Neurocomputing* 37(2001)71-90.
- [9] A S Weigend, D E Rumelhart and B A Huberman, Generalization by weight-elimination applied to currency exchange rate prediction. *Proc. Intl Joint Conf. on Neural Networks* 1(Seattle, 19916) 837-841.
- [10] Y Shin and J Ghosh, Approximation of multivariate functions using ridge polynomial networks, *International Joint Conference on Neural Networks* 2 (1992) 380-385.
- [11] M Sinha, K Kumar and P K Kalra, Some new neural network architectures with improved learning schemes. *Soft Computing*, 4 (2000) 214-223.
- [12] R Setiono, A penalty-function approach for pruning feed forward neural networks, *Neural Networks* 9 (1997) 185-204.
- [13] W Wu and Y S Xu, Deterministic convergence of an online gradient method for neural networks, *Journal of Computational and Applied Mathematics* 144 (1-2) (2002) 335-347.
- [14] H S Zhang and W Wu, Boundedness and convergence of online gradient method with penalty for linear output feed forward neural networks, *Neural Process Letters* 29 (2009) 205-212.
- [15] H F Lu, W Wu, C Zhang and X Yan, Convergence of Gradient Descent Algorithm for Pi- Sigma Neural Networks, *Journal of Information and Computational Science* 3: 3 (2006) 503-509.
- [16] YX Yuan and WY Sun, *Optimization Theory and Methods*, Science Press, Beijing, 2001.
- [17] J Kong and W Wu, Online gradient methods with a punishing term for neural networks. *Northeast Math. J* 1736(2001) 371-378.
- [18] W Wu, G R Feng and X Z Li, Training multiple perceptrons via minimization of sum of ridge functions, *Advances in Computational Mathematics* 17 (2002) 331-347.