

行政院國家科學委員會專題研究計畫 成果報告

相關性隱藏節點與學習演算法 研究成果報告(精簡版)

計畫類別：個別型
計畫編號：NSC 95-2416-H-004-049-
執行期間：95年08月01日至96年07月31日
執行單位：國立政治大學資訊管理學系

計畫主持人：蔡瑞煌

計畫參與人員：學士級-專任助理：楊佳鳳

報告附件：出席國際會議研究心得報告及發表論文

處理方式：本計畫可公開查詢

中華民國 96 年 07 月 30 日

1. The Construction Procedure

Huang and Babri (1998) propose an elegant construction method to set up a real-valued single-hidden layer feed-forward neural network (SLFN) with N hidden nodes that successfully learns N distinct samples with zero error. In a *correlated* real-valued single-hidden layer feed-forward neural network (SLFN), the weight vectors in the input layer of all its hidden nodes are linearly dependent. Tsaih and Wan (2007) realize that the SLFN constructed by Huang and Babri (1998) is correlated. They further show that the correlated SLFN has the property of hyperplane preimages. The correlated SLFN provides a hyperplane-preimage approach for the nonlinear regression problem with the assumption of linear preimage. Such usages motivate a study of the construction procedure for creating a correlated SLFN with less than N hidden nodes that perfectly fits N distinct samples.

The proposed construction procedure will initially set up one hidden node and then recruit (add) more (linearly dependent) hidden nodes during the learning process. In the literature, there are some similar procedures; for instance, the tiling algorithm for binary-valued layered feed-forward neural networks (cf. Me'zard and Nadal, 1989), the cascade-correlation algorithm (cf. Fahlman and C. Lebiere, 1990), and the upstart algorithm for binary-valued layered feed-forward neural networks (cf. Frean, 1990). In contrast to these researches, this study copes with the correlated real-valued SLFN.

In the context of *estimation*, the response y equates $f(\mathbf{x}, \mathbf{w}) + \delta$ where \mathbf{w} is the parameter vector and δ is the error term. Usually, the function form of f is predetermined and fixed during the process of deriving its associated \mathbf{w} from a given data set of observations $\{(1\mathbf{x}, 1y), \dots, (N\mathbf{x}, Ny)\}$, with c_y being the observed response corresponding to the c^{th} observation $c\mathbf{x}$.

The least squares estimator (LSE) is one of the most popular methods for estimating. If $\hat{\mathbf{w}}$ denotes any estimate of \mathbf{w} , then LSE is defined to minimize $\sum_{c=1}^N c e^2$, where

$$c e = c y - f(c\mathbf{x}, \hat{\mathbf{w}}). \quad (1)$$

The generalized delta rule proposed in (Rumelhart, Hilton, and Williams, 1986) for the learning process of SLFN is a kind of (nonlinear) LSE. The LSE, however, is known to be very sensitive to outliers.

In the literature of linear regression analysis, there are two approaches of dealing with outlier problems: deletion diagnostics and robust estimators (cf. Rousseeuw and Leroy (1987, page 8)). The diagnostic approach assesses the influence of an individual observation or a subset of observations to the LSE. The diagnostic approach is useful to assess the adequacy of the underlying assumption and to identify unexpected characteristics of the data. One way for the diagnostics is to identify the observations that leave the largest change in the diagnostic quantity (cf. (Cook and Weisberg, 1982)(Atkinson, 1985)) when they are excluded from the fitted data set.

As for the robust statistics approach, the robustness analysis (cf. (Hampel, 1986)) limits the attention to a "trimmed" sum of squared residuals instead of adding all the squared residuals as in the LSE. If only the first q of those ordered squared residuals are included in the summation, then the least trimmed squares (LTS) estimator is defined as

$$\text{Minimize } \sum_{c=1}^q [c] e^2, \quad (2)$$

where $_{[c]}e^2$ denotes the ordered squared residuals; that is, $_{[1]}e^2 \leq _{[2]}e^2 \leq \dots \leq _{[N]}e^2$. Zaman, Rousseeuw, and Orthan (2001) suggest that $\lfloor 0.75N \rfloor^1$ is a reasonable value for q in most empirical studies.

Atkinson and Cheng (1999) adapt the forward search algorithm proposed in (Atkinson, 1994) to develop the LTS estimates. The forward search algorithm consists of randomly adopting an (initial) subset of $m+1$ observations to fit the linear regression model, ordering the residuals of all N observations, and then augmenting the subset gradually by including extra observations based upon the smallest squared residuals principle.

The C-step² of Rousseeuw and Van Driessen (2002) can release quite fast a series of subsets of observations whose corresponding total squared residuals are refined gradually. The last subset results in a good linear fitting function which is an approximation of the LTS estimator.

2. The Mapping Requirement and Notations

An SLFN provides a nonlinear mapping between \mathbf{x} and y , whose form is $y = f(\mathbf{x})$. f is a nonlinear function whose parameters (i.e., weights and biases) are derived from a given data set of mapping samples $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$ with $\mathbf{x}^{c_1} \neq \mathbf{x}^{c_2}$, $c_1 \neq c_2$, and with t^c the target value of y corresponding to \mathbf{x}^c . $\mathbf{x} \equiv (x_1, x_2, \dots, x_m)^T \in R^m$ where x_j is the j^{th} input component, with j from 1 to m .

Hereafter, m and p denote the numbers of adopted input and hidden nodes, respectively; $w_{j_0}^2$ stands for the bias of the j^{th} hidden node; $\mathbf{w}_j^2 \equiv (w_{j_1}^2, w_{j_2}^2, \dots, w_{j_m}^2)^T$ for the weights between the j^{th} hidden node and input layer; w_0^3 for the bias of the output node; and $\mathbf{w}^3 \equiv (w_1^3, w_2^3, \dots, w_p^3)^T$ for the weights between the output node and all hidden nodes. Characters in bold represent column vectors; the superscript T indicates transposition.

Let the $\tanh(t)$ activation function be used by all hidden nodes and a linear activation function be used by the output node. Thus, given the c^{th} sample \mathbf{x}^c , the activation value of the j^{th} hidden node $a^c(w_{j_0}^2, \mathbf{w}_j^2)$ and the output value y^c are as follows:

$$a^c(w_{j_0}^2, \mathbf{w}_j^2) \equiv \tanh(w_{j_0}^2 + \sum_{i=1}^m w_{ji}^2 x_i^c); \quad (3)$$

$$y^c \equiv w_0^3 + \sum_{j=1}^p w_j^3 a^c(w_{j_0}^2, \mathbf{w}_j^2). \quad (4)$$

Given the N mapping samples, let $\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \equiv (a^1(w_{j_0}^2, \mathbf{w}_j^2), a^2(w_{j_0}^2, \mathbf{w}_j^2), \dots, a^N(w_{j_0}^2, \mathbf{w}_j^2))^T \in (-1,1)^N$ be the responding vector of the j^{th} hidden node with the c^{th} component being $a^c(w_{j_0}^2, \mathbf{w}_j^2)$. Furthermore, let $\mathbf{1}$ be a $N \times 1$ vector with all components 1 and $\mathbf{T} \equiv (t^1, t^2, \dots, t^N)^T \in R^N$. Thus, the set of simultaneous equations $w_0^3 + \sum_{j=1}^p w_j^3 \tanh(w_{j_0}^2 + \sum_{i=1}^m w_{ji}^2 x_i^c) = t^c \forall c = 1, \dots, N$ is equivalent to system (5), which states that \mathbf{T} is in the space spanned by $\mathbf{1}$ and p responding vectors, $\{\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), j = 1, \dots, p\}$.

$$w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) = \mathbf{T}. \quad (5)$$

¹ Hereafter, $\lfloor x \rfloor$ is the largest integer not larger than x .

² C stands for "concentration". The idea of C-step has been implemented in the built-in function *lts.reg* of Splus which is a commercial statistical computing package published by MathSoft Co..

Hereafter, let ${}_a\mathbf{R}_b \equiv \mathbf{a} - \frac{\mathbf{b}^T \mathbf{a}}{\|\mathbf{b}\|} \frac{\mathbf{b}}{\|\mathbf{b}\|}$ denote the residual of vector \mathbf{a} regarding \mathbf{b} after the part parallel with the vector \mathbf{b} has been taken away. $\frac{\mathbf{b}^T \mathbf{a}}{\|\mathbf{b}\|} \frac{\mathbf{b}}{\|\mathbf{b}\|}$ is the projection of \mathbf{a} in the direction of \mathbf{b} ; and ${}_a\mathbf{R}_b$ is orthogonal to \mathbf{b} since $\mathbf{b}^T {}_a\mathbf{R}_b = \mathbf{b}^T (\mathbf{a} - \frac{\mathbf{b}^T \mathbf{a}}{\|\mathbf{b}\|} \frac{\mathbf{b}}{\|\mathbf{b}\|}) = 0$.

Similarly, let ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} \equiv {}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} - \frac{({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|}$ denote the residual of vector \mathbf{a} regarding the ordered sequence of linearly independent vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$ after the part in the (sub-)space spanned by $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$ has been taken away. $\frac{({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|}$ is the projection of ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}$ in the direction of ${}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}$; and ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$ is orthogonal to ${}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}$ since $({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}})^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} = ({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} - \frac{({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|})^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} = 0$. By definition we have

Lemma 1. Furthermore, $\mathbf{a} = \frac{\mathbf{a}^T \mathbf{b}^1}{\|\mathbf{b}^1\|} \frac{\mathbf{b}^1}{\|\mathbf{b}^1\|} + \sum_{j=2}^k \frac{({}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}})^T {}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} \frac{{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} + {}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$

and **Lemma 2** lists some properties associated with the above proposed residual vector.

Lemma 1: If \mathbf{a} is linearly dependent with \mathbf{b} , then ${}_a\mathbf{R}_b = \mathbf{0}$. Similarly, if \mathbf{a} can be linearly represented by the set of vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$,³ then ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} = \mathbf{0}$.

Lemma 2: (i) ${}_a\mathbf{R}_b$ is orthogonal to \mathbf{b} . (ii) ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$ is orthogonal to the subspace spanned by the set of linearly independent vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$. (iii) If ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} = \mathbf{0}$, then ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k+1}\}} = \mathbf{0}$.

3. The Proposed Construction Procedure

Table 1 presents the proposed procedure for constructing a correlated SLFN appropriate for fitting the mapping embedded in $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$.

Table 1. The proposed deterministic procedure for constructing an appropriate correlated SLFN for the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$. $\mathbf{T} \equiv (t^1, t^2, \dots, t^N)^T$ and $\mathbf{v}(w_0, \mathbf{w}) \equiv (a^1(w_0, \mathbf{w}), a^2(w_0, \mathbf{w}), \dots, a^N(w_0, \mathbf{w}))^T$ is a $N \times 1$ vector with the c^{th} component being $a^c(w_0, \mathbf{w}) \equiv \tanh(\mathbf{w}^T \mathbf{x}^c + w_0)$, in which $\mathbf{w} \equiv (w_1, w_2, \dots, w_m)^T$.

Step 1: Calculate ${}_T\mathbf{R}_1$. If ${}_T\mathbf{R}_1 = \mathbf{0}$, then (i) claim that the fitting job requests no hidden node; (ii) set the bias of output node as $\frac{\mathbf{T}^T \mathbf{1}}{\mathbf{1}^T \mathbf{1}}$ and the weight vector between the output

³ Namely, \mathbf{a} is in the (sub-)space spanned by $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$.

node and the input layer as $\mathbf{0}$; and (iii) stop.

Step 2: Apply the C-step to all N observations to obtain the $m+1$ input samples that are linearly independent. Let $\mathbf{I}(m+1)$ be the set of indices of these samples and $\mathbf{I}(N)$ be the set of indices of all samples.

Step 3: Calculate \tilde{t}^c which equates $\tanh^{-1}\left(\frac{t^c - \min_{c \in \mathbf{I}(N)} t^c + 1}{\max_{c \in \mathbf{I}(N)} t^c - \min_{c \in \mathbf{I}(N)} t^c + 2}\right)$ from $\{t^c: \forall c \in \mathbf{I}(m+1)\}$.

Next, apply the linear regression method to the data set $\{(\mathbf{x}^c, \tilde{t}^c): \forall c \in \mathbf{I}(m+1)\}$ to get a set of $m+1$ weights.

Step 4: Set one hidden node in the network whose values of w_{10}^2 and w_1^2 are assigned as the values of the weights obtained in Step 3, and initial values of w_0^3 and w_1^3 are assigned as $\min_{c \in \mathbf{I}(N)} t^c - 1$ and $\max_{c \in \mathbf{I}(N)} t^c - \min_{c \in \mathbf{I}(N)} t^c + 2$, respectively. Then set $\gamma_1 = 1$ and $p = 1$.

Step 5: If $w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) = \mathbf{T}$, then (i) claim that the fitting job requests p hidden nodes with the bias and weights being the above w_0^3 , w_{j0}^2 , \mathbf{w}_j^2 , and w_j^3 for all $1 \leq j \leq p$; and (ii) stop.

Step 6: If $w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) \neq \mathbf{T}$, then solve $\min_{w_0, \gamma} \|\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}\|^2$ and let $(w_0^*, \gamma^*) \equiv \arg(\min_{w_0, \gamma} \|\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}\|^2)$.

Step 7: set $\gamma_{p+1} = \gamma^*$, $w_{p+1,0}^2 = w_0^*$, $\mathbf{w}_{p+1}^2 = \gamma_{p+1} \mathbf{w}_1^2$, and $p+1 \rightarrow p$.

Step 8: Apply the linear regression method to the data set $\{((a^c(w_{10}^2, \mathbf{w}_1^2), a^c(w_{20}^2, \mathbf{w}_2^2), \dots, a^c(w_{p0}^2, \mathbf{w}_p^2))^T, t^c): \forall c \in \mathbf{I}(N)\}$ to get values of w_0^3 and $w_j^3, \forall j = 1, \dots, p$. Then go to Step 5.

Step 2 releases $m+1$ linearly independent input samples through applying the C-step to all N input samples.⁴ Let $\mathbf{I}(m+1)$ be the set of indices of these (linearly independent) input

samples. Step 3 calculates \tilde{t}^c via $\tanh^{-1}\left(\frac{t^c - \min_{c \in \mathbf{I}(N)} t^c + 1}{\max_{c \in \mathbf{I}(N)} t^c - \min_{c \in \mathbf{I}(N)} t^c + 2}\right) \forall c \in \mathbf{I}(m+1)$. Then Step 3

applies the linear regression method to the data set $\{(\mathbf{x}^c, \tilde{t}^c): \forall c \in \mathbf{I}(m+1)\}$ to get the unique solution of $(w_{10}^2, \mathbf{w}_1^2)$ of system (6), which is a system of $m+1$ linear equations in $m+1$ unknowns.

⁴ The choices of the subset at Step 2 can be adapted by other considerations (cf. Stromberg (1993)).

$$w_{10}^2 + \sum_{j=1}^m w_{1j}^2 x_j^c = \tilde{t}^c \quad \forall c \in \mathbf{I}(m+1). \quad (6)$$

Step 4 sets up the network with one hidden node whose values of $(w_{10}^2, \mathbf{w}_1^2)$ are assigned as obtained in Step 3. The initial values of w_0^3 and w_1^3 are assigned as $\min_{c \in \mathbf{I}(N)} t^c - 1$ and $\max_{c \in \mathbf{I}(N)} t^c - \min_{c \in \mathbf{I}(N)} t^c + 2$, respectively. According to Rousseeuw and Van Driessen (2002), this setup network renders ${}^c e^2 = 0 \quad \forall c \in \mathbf{I}(m+1)$ and is a good approximation of the LTS estimator.

Step 5 denotes the stopping criterion of the proposed procedure.

The minimization in Step 6 and the assignment in Step 7 determine the bias and weights for the connections of the input nodes to the most newly recruited hidden node. All biases and weights for the connections of the input nodes to the previously recruited hidden nodes are unchanged. Furthermore, the assignment of Step 7 renders the constructed SLFN correlated since $\mathbf{w}_j^2 = \gamma_j \mathbf{w}_1^2, j=1, \dots, p$.

4. The Correctness of the Proposed Procedure

We now prove that the correlated SLFN constructed by the procedure stated in Table 1 meets the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$ without error.

Tsaih and Wan (2007) state that, for any given $\{(\mathbf{x}^c, t^c): \forall c = 1, \dots, N\}$, there exists a set of $\{(w_{j0}^2, \gamma_j \mathbf{w}_1^2), j=1, \dots, N-1\}$ such that the associated square matrix $(\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{N-1,0}^2, \gamma_{N-1} \mathbf{w}_1^2))$ is invertible and the mapping requirement is achieved perfectly. Therefore, we have **Lemma 3** and the proposed procedure will stop at any p with $0 \leq p \leq N-1$.

Lemma 3: If $p < N-1$, then there always exist w_0 and γ such that ${}_{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} \neq \mathbf{0}$.

Proof of Lemma 3: Suppose there is a $p < N-1$ such that there are no w_0 and γ to render ${}_{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} \neq \mathbf{0}$. In other words, $p < N-1$ and there are no w_0 and γ such that $\mathbf{v}(w_0, \gamma \mathbf{w}_1^2)$ is linearly independent with $\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \mathbf{v}(w_{20}^2, \mathbf{w}_2^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}$. This contradicts with the statement of Tsaih and Wan (2007). **Q.E.D.**

When the procedure stops at Step 1, the SLFN constructed at Step 1 meets the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$ without error since ${}_{\mathbf{T}} \mathbf{R}_1 = \mathbf{0}$. On the other hand, it is obvious to have **Lemma 4**, which states the necessary condition of a SLFN with p hidden nodes appropriate for the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$. Thus the stopping criterion stated in Step 5 is suitable.

Lemma 4: Regarding the mapping samples of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$, the SLFN with p hidden nodes is appropriate for the mapping requirement if $w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) = \mathbf{T}$.

Consider the case of $w_0^3 \mathbf{1} + \sum_{j=1}^{p-1} w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) \neq \mathbf{T}$ and $w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) = \mathbf{T}$. Namely, ${}_{\mathbf{T}} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)\}} \neq \mathbf{0}$ for each $1 \leq j \leq p-1$ and ${}_{\mathbf{T}} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} = \mathbf{0}$. For each $1 \leq j \leq p-1$, from calculation, we have ${}_{\mathbf{T}} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0^2, \gamma \mathbf{w}_1^2)\}} = {}_{\mathbf{T}} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)\}} - \frac{({}_{\mathbf{T}} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)\}})^T {}_{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)\}}}{\|{}_{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)\}}\|} \quad \text{and}$

$$\begin{aligned} & \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} \right\|^2 = \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|^2 \left(1 - \left(\frac{(\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}})^{\mathbf{T}}}{\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|} \right. \right. \\ & \left. \left. \frac{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}}{\left\| \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|} \right)^2 \right). \quad \text{Hence we have Lemma 5. Furthermore, suppose} \\ & \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}} \neq \mathbf{0}, \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2)\}} \right\|^2 = 0 \text{ if and only if } \left(\frac{(\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}})^{\mathbf{T}}}{\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}} \right\|} \right. \\ & \left. \frac{\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}}}{\left\| \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}} \right\|} \right)^2 = 1, \text{ which implies the following Lemma 6 is true.} \end{aligned}$$

Lemma 5: If $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \neq \mathbf{0}$, then $\min_{w_0, \gamma} \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} \right\|^2 > \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|^2$.

Proof of Lemma 5: From Lemma 3, there always exists w_0 and γ such that $\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}$ is a non-zero vector. Now $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} = \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} - \frac{(\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}})^{\mathbf{T}} \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}}{\left\| \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|} \frac{\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}}{\left\| \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|}$ and $\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}$ and $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}$ are orthogonal to each other. Thus, $\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|^2 > \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} \right\|^2$. In other words, the optimization of $\min_{w_0, \gamma} \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} \right\|^2$ leads to a non-zero $\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}}$, in which $(w_0^*, \gamma^*) \equiv \arg(\min_{w_0, \gamma} \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}} \right\|^2)$ and $\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2), \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2)\}} \right\|^2 < \left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)\}} \right\|^2$.

Q.E.D.

Lemma 6: Suppose $\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}} \right\|^2 \neq 0$. $\left\| \mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2)\}} \right\|^2 = 0$ if and only if $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}}$ and $\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_0^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}}$ are parallel.

Lemma 7: The (sub-)space spanned by the set of linearly independent vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$ is equivalent with the one spanned by the set of orthonormal vectors $\left\{ \frac{\mathbf{b}^1}{\|\mathbf{b}^1\|}, \frac{\mathbf{b}^2 \mathbf{R}_{\mathbf{b}^1}}{\|\mathbf{b}^2 \mathbf{R}_{\mathbf{b}^1}\|}, \dots, \frac{\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \right\}$.

Proof of Lemma 7: Let us prove by induction. It is trivial for the case of any set of two linearly independent vectors $\{\mathbf{b}^1, \mathbf{b}^2\}$ since $\mathbf{b}^2 \mathbf{R}_{\mathbf{b}^1}$ is orthogonal to \mathbf{b}^1 and $\mathbf{b}^1 \mathbf{R}_{\mathbf{b}^2}$ is orthogonal to \mathbf{b}^2 .

Now consider the case of any set of $k+1$ linearly independent vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^{k+1}\}$ with $k \geq 2$. Assume the subspace spanned by the subset of vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$ is equivalent with the one spanned by the set of orthonormal vectors $\left\{ \frac{\mathbf{b}^1}{\|\mathbf{b}^1\|}, \frac{\mathbf{b}^2 \mathbf{R}_{\mathbf{b}^1}}{\|\mathbf{b}^2 \mathbf{R}_{\mathbf{b}^1}\|}, \dots, \frac{\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \right\}$. The

vector \mathbf{b}^{k+1} can be represented as $\frac{(\mathbf{b}^{k+1})^T \mathbf{b}^1}{\|\mathbf{b}^1\|} \frac{\mathbf{b}^1}{\|\mathbf{b}^1\|} + \sum_{j=2}^k \frac{(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}})^T \mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|\mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} \frac{\mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|\mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} + \mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$, in which $\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$ is a non-zero vector since \mathbf{b}^{k+1} is linearly independent with the set of vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$. Thus the (sub-)space spanned by the set of vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^{k+1}\}$ is equivalent with the one spanned by the set of vectors $\{\frac{\mathbf{b}^1}{\|\mathbf{b}^1\|}, \frac{\mathbf{b}^2 \mathbf{R}_{\{\mathbf{b}^1\}}}{\|\mathbf{b}^2 \mathbf{R}_{\{\mathbf{b}^1\}}\|}, \dots, \mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}\}$. Furthermore, $(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}})^T \mathbf{b}^1 = (\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} - \frac{(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T \mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|})^T \mathbf{b}^1 = 0$; and, for any $2 \leq j \leq k-1$, $(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}})^T \mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}} = (\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} - \frac{(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T \mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|})^T \mathbf{b}^j \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}} = 0$; and $(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}})^T \mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} = (\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} - \frac{(\mathbf{b}^{k+1} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}})^T \mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|\mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|})^T \mathbf{b}^k \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}} = 0$. Q.E.D.

Let $\mathbf{u}^0 \equiv \frac{\mathbf{1}}{\|\mathbf{1}\|}$, $\mathbf{u}^1 \equiv \frac{\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \mathbf{R}_1}{\|\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \mathbf{R}_1\|}$, and $\mathbf{u}^j \equiv \frac{\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{j_0}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j-1,0}^2, \mathbf{w}_{j-1}^2)\}}}{\|\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{j_0}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{j-1,0}^2, \mathbf{w}_{j-1}^2)\}}\|}$ for all $2 \leq j \leq p$.

Thus $\mathbf{T}^T \mathbf{u}^k = (w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2))^T \mathbf{u}^k = \sum_{j=k}^p w_j^3 (\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2))^T \mathbf{u}^k$ for all $1 \leq k \leq p$, since, from Lemma 7, $\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2)$ is in the subspace spanned by $\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^j\}$ and $(\mathbf{u}^j)^T \mathbf{u}^k \neq 0 \forall k \leq j \leq p$.

Thus $w_p^3 \equiv \frac{\mathbf{T}^T \mathbf{u}^p}{(\mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2))^T \mathbf{u}^p}$, $w_j^3 \equiv \frac{(\mathbf{T} - \sum_{l=j+1}^p w_l^3 \mathbf{v}(w_{l_0}^2, \mathbf{w}_l^2))^T \mathbf{u}^j}{(\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2))^T \mathbf{u}^j}$ for all $p-1 \geq j \geq 1$, and $w_0^3 \equiv \frac{(\mathbf{T} - \sum_{l=1}^p w_l^3 \mathbf{v}(w_{l_0}^2, \mathbf{w}_l^2))^T \mathbf{u}^0}{\mathbf{1}^T \mathbf{u}^0}$.

Let $\mathbf{A}^{-1} \equiv \mathbf{I}$ and $\mathbf{A}^k \equiv \mathbf{A}^{k-1} - \mathbf{u}^k (\mathbf{u}^k)^T$ for all $0 \leq k \leq N-1$. Thus, $\mathbf{A}^k = \mathbf{I} - \sum_{j=0}^k \mathbf{u}^j (\mathbf{u}^j)^T$ for all $0 \leq k \leq N-1$. Since all vectors in the set $\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^k\}$ are orthonormal, $(\mathbf{A}^k)^T = \mathbf{A}^k$ and $\mathbf{A}^k \mathbf{A}^k = \mathbf{A}^k$. Furthermore, because $\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{T} - \sum_{j=0}^p (\mathbf{u}^j)^T \mathbf{T} \mathbf{u}^j$, $\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{A}^p \mathbf{T}$. Similarly, $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{A}^p \mathbf{v}(w_0, \mathbf{w})$. Thus, $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}^T \mathbf{T} = \mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{T}$, $\mathbf{v}(w_0, \mathbf{w})^T \mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{T}$, $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}^T \mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{T}$, $\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}^T \mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{T}^T \mathbf{A}^p \mathbf{T}$, and $\|\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}\|^2 = \mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{v}(w_0, \mathbf{w})$. thus, $\|\mathbf{T} \mathbf{R}_{\mathbf{v}(w_0, \mathbf{w})}\|^2 = \mathbf{T}^T \mathbf{T} - \frac{(\mathbf{v}(w_0, \mathbf{w})^T \mathbf{T})^2}{\mathbf{v}(w_0, \mathbf{w})^T \mathbf{v}(w_0, \mathbf{w})}$ and $\|\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{v}(w_0, \mathbf{w})\}}\|^2 = \mathbf{T}^T \mathbf{A}^p \mathbf{T} - \frac{(\mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{T})^2}{\mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{v}(w_0, \mathbf{w})} = \frac{\mathbf{T}^T \mathbf{A}^p \mathbf{T} \mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{v}(w_0, \mathbf{w}) - (\mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{T})^2}{\mathbf{v}(w_0, \mathbf{w})^T \mathbf{A}^p \mathbf{v}(w_0, \mathbf{w})}$.

Lemma 8: If $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{j_0}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2)\}} \neq \mathbf{0}$ and there exists a vector $\mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}_1^2)$ such that $\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{j_0}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2), \mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}_1^2)\}} = \mathbf{0}$, then $\min_{w_0, \mathcal{J}} \|\mathbf{T} \mathbf{R}_{\{\mathbf{1}, \mathbf{v}(w_{j_0}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p_0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0, \mathcal{J} \mathbf{w}_1^2)\}}\|^2$ leads to a non-zero

$\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}$ such that $\frac{\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}}{\|\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}\|}$ and $\frac{\mathbf{v}(\bar{w}_0, \bar{\gamma} \bar{\mathbf{w}}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}}{\|\mathbf{v}(\bar{w}_0, \bar{\gamma} \bar{\mathbf{w}}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}\|}$ are parallel and $\|\mathbf{T} \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2)\}}\|^2 = 0$.

Proof: Let $\bar{\mathbf{u}}^{p+1} \equiv \frac{\mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$. Thus, facts of $\mathbf{T} \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} \neq \mathbf{0}$ and $\mathbf{T} \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(\bar{w}_0, \bar{\gamma} \bar{\mathbf{w}}_1^2)\}} = \mathbf{0}$ imply that $\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}$ is an orthonormal set, $\mathbf{T} = \sum_{j=0}^p (\mathbf{u}^j)^T \mathbf{T} \mathbf{u}^j + (\bar{\mathbf{u}}^{p+1})^T \mathbf{T} \bar{\mathbf{u}}^{p+1}$, and $(\bar{\mathbf{u}}^{p+1})^T \mathbf{T} \neq 0$.

$$\text{Suppose } \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \equiv \sum_{j=0}^p (\mathbf{u}^j)^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{u}^j + (\bar{\mathbf{u}}^{p+1})^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \bar{\mathbf{u}}^{p+1} + \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}}.$$

$$\text{Thus, } \|\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}\|^2 = \frac{\mathbf{T}^T \mathbf{A}^p \mathbf{T} (\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) - ((\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{T})^2}{(\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} =$$

$$\frac{1}{(\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} [(\mathbf{T}^T \bar{\mathbf{u}}^{p+1})^2 (\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T ((\bar{\mathbf{u}}^{p+1})^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \bar{\mathbf{u}}^{p+1} + \mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}})$$

$$- ((\bar{\mathbf{u}}^{p+1})^T \mathbf{T} (\bar{\mathbf{u}}^{p+1})^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^2] = \frac{(\mathbf{T}^T \bar{\mathbf{u}}^{p+1})^2}{(\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} [(((\bar{\mathbf{u}}^{p+1})^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^2 +$$

$$\|\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}}\|^2) - ((\bar{\mathbf{u}}^{p+1})^T \mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^2] = \frac{(\mathbf{T}^T \bar{\mathbf{u}}^{p+1})^2}{(\mathbf{v}(w_0, \gamma \mathbf{w}_1^2))^T \mathbf{A}^p \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)} \|\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}}\|^2.$$

If $\|\mathbf{v}(w_0, \gamma \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}}\|^2 > 0$, then $\|\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}\|^2 > \|\mathbf{T} \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{v}(\bar{w}_0, \bar{\gamma} \bar{\mathbf{w}}_1^2)\}}\|^2$ because of a non-zero $(\bar{\mathbf{u}}^{p+1})^T \mathbf{T}$. Therefore, the optimization of $\min_{w_0, \gamma} \|\mathbf{T} \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0, \gamma \mathbf{w}_1^2)\}}\|^2$ leads to a non-zero $\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}}$ such that $\|\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p, \bar{\mathbf{u}}^{p+1}\}}\|^2 = 0$. From **Lemma 6**, $\frac{\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$ and $\frac{\mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{v}(\bar{w}_0, \bar{\mathbf{w}}) \mathbf{R}_{\{\mathbf{u}^0, \mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$ are parallel and thus, $\|\mathbf{T} \mathbf{R}_{\{1, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{v}(w_0^*, \gamma^* \mathbf{w}_1^2)\}}\|^2 = 0$. **Q.E.D.**

References

1. M. Me'zard and J. Nadal, *Learning in feedforward layered networks: The tiling algorithm. Journal of Physics. A* **22**, 2191– 2204 (1989).
2. S. Fahlman and C. Lebiere, *The Cascade-Correlation Learning Architecture. In Touretzky, D. (Eds.), Advances in Neural Information Processing Systems II (Denver, 1989) Morgan Kaufmann, San Mateo* (1990).
3. M. Frenn, *The Upstart Algorithm: A Method for Constructing and Training Feedforward Neural Networks. Neural Computation.* **2**, 198– 209(1990).
4. G. Huang and H. Babri, *Upper bounds on the number of hidden neurons in feedforward networks with arbitrary bounded nonlinear activation functions. IEEE Transactions on Neural Networks.* **9**, 224– 229(1998).
5. D. Rumelhart, G. Hinton, and R. Williams, "Modeling Internal Representations By Error Propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of*

- Cognition*, vol. 1, D. Rumelhart and J. McClelland, Eds. Cambridge, MA: MIT Press, 1986, pp. 318-362.
6. P. J. Rousseeuw and A. M. Leroy, *Robust Regression and Outlier Detection*. New York: Wiley, 1987.
 7. R. D. Cook and S. Weisberg, *Residuals and Influence in Regression*, London: Chapman and Hall, 1982.
 8. A. C. Atkinson, *Plots, Transformations and Regression*, Oxford: Oxford University Press, 1985.
 9. F. R. Hampel, E. M. Ronchetti, P. J. Rousseeuw and W. A. Stahel, *Robust Statistics: The Approach Based on Influence Functions*. New York: John Wiley, 1986.
 10. A. Zaman, P. J. Rousseeuw and M. Orhan, "Econometric Applications of High-Breakdown Robust Regression Techniques," *Econometrics Letters*, vol. 71, pp. 1-8, Apr. 2001.
 11. A. C. Atkinson and T. C. Cheng, "Computing Least Trimmed Squares Regression with the Forward Search," *Statistics and Computing*, vol. 9, pp. 251-263, Nov. 1999.
 12. A. C. Atkinson, "Fast Very Robust Methods for the Detection of Multiple Outliers," *Journal of the American Statistical Association*, vol. 89, pp. 1329-1339, 1994.
 13. P. J. Rousseeuw and K. Van Driessen, "Computing LTS Regression for Large Data Sets," *Estadistica*, vol. 54, 163-190, 2002.
 14. A. J. Stromberg, "Computation of High Breakdown Nonlinear Regression Parameters," *Journal of the American Statistical Association*, vol. 88, pp. 237-244, Mar. 1993.

計畫成果自評：

此研究計畫成果豐碩，已被送到相關研討會發表。其延伸之研究亦在進行中。

出席國際學術會議心得報告

計畫編號	95-2416-H-004-049
計畫名稱	相關性隱藏節點與學習演算法
出國人員姓名 服務機關及職稱	蔡瑞煌 國立政治大學資訊管理學系, 教授
會議時間地點	July 18-24, 2007, Salt Lake City, U.S.A.
會議名稱	Joint Conference on Information Sciences
發表論文題目	A Constructive Learning Procedure

一、參加會議經過

我於 7/19 晚上到達 Salt Lake City。我於 7/20 主持 section CIEF-III 並於其間發表論文。附件一是相關議程。

我也於 7/20~7/21 聆聽此次會議裡的 Keynote Speeches 和論文發表。

二、與會心得

Keynote Speeches 邀請了不少的資訊科學學術界裡之知名學者來演講，我受益不少。

A Constructive Learning Procedure*

RAY TSAIH

Department of Management Information Systems, National Chengchi University, No.64, Sec. 2, Jhihnan Rd., Wunshan District, Taipei
City 116, Taiwan

This study explores a deterministic learning procedure for the realization of a real-valued single-hidden layer feed-forward neural network (SLFN) with tanh activation functions of the hidden-layer nodes for arbitrary mapping problems.

1. The Constructive Learning Procedure

The proposed learning procedure will use none hidden node initially and recruit (add) more hidden nodes during the learning process. The goal of the proposed constructive learning procedure is to create a SLFN for fitting perfectly all given mapping samples. In the literature, there are some similar procedures; for instance, the tiling algorithm for binary-valued layered feed-forward neural networks (cf. [1]), the cascade-correlation algorithm (cf. [2]), and the upstart algorithm for binary-valued layered feed-forward neural networks (cf. [3]). In contrast to these researches, this study copes with the real-valued SLFN.

2. The Mapping Requirement and Notations

An SLFN provides a nonlinear mapping between \mathbf{x} and y , whose form is $y = f(\mathbf{x})$. f is a nonlinear function whose parameters (i.e., weights and biases) are derived from a given data set of mapping samples $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$ with $\mathbf{x}^{c1} \neq \mathbf{x}^{c2}$, $c_1 \neq c_2$, and with t^c the target value of y corresponding to \mathbf{x}^c . $\mathbf{x} \equiv (x_1, x_2, \dots, x_m)^T \in R^m$ where x_j is the j^{th} input component, with j from 1 to m .

Hereafter, m and p denote the numbers of adopted input and hidden nodes, respectively; $w_{j_0}^2$ stands for the bias of the j^{th} hidden node; $\mathbf{w}_j^2 \equiv (w_{j_1}^2, w_{j_2}^2, \dots, w_{j_m}^2)^T$ for the weights between the j^{th} hidden node and input layer; w_0^3 for the bias of the output node; and $\mathbf{w}^3 \equiv (w_1^3, w_2^3, \dots, w_p^3)^T$ for the weights between the output node and all hidden nodes. Characters in bold represent column vectors; the superscript T indicates transposition.

Let the $\tanh(t)$ activation function be used by all hidden nodes and a linear activation function be used by the output node. Thus, given the c^{th} sample \mathbf{x}^c , the activation value of the j^{th} hidden node $a^c(w_{j_0}^2, \mathbf{w}_j^2)$ and the output value y^c are as follows:

$$a^c(w_{j_0}^2, \mathbf{w}_j^2) \equiv \tanh(w_{j_0}^2 + \sum_{i=1}^m w_{ji}^2 x_i^c); \quad (1)$$

$$y^c \equiv w_0^3 + \sum_{j=1}^p w_j^3 a^c(w_{j_0}^2, \mathbf{w}_j^2). \quad (2)$$

Given the N mapping samples, let $\mathbf{v}(w_{j_0}^2, \mathbf{w}_j^2) \equiv (a^1(w_{j_0}^2, \mathbf{w}_j^2), a^2(w_{j_0}^2, \mathbf{w}_j^2), \dots, a^N(w_{j_0}^2, \mathbf{w}_j^2))^T \in (-1, 1)^N$ be the responding vector of the j^{th} hidden node with the c^{th} component being $a^c(w_{j_0}^2, \mathbf{w}_j^2)$. Furthermore, let $\mathbf{1}$ be a $N \times 1$ vector with all components 1 and $\mathbf{T} \equiv (t^1, t^2, \dots, t^N)^T \in R^N$. Thus, the set of simultaneous equations $w_0^3 + \sum_{j=1}^p w_j^3 \tanh(w_{j_0}^2 + \sum_{i=1}^m w_{ji}^2 x_i^c) = t^c \quad \forall c = 1, \dots, N$ is equivalent to system (3), which states that \mathbf{T} is in the (sub-)space spanned by $\mathbf{1}$ and p responding vectors,

* This study is supported by National Science Council of R.O.C. under Grants No. NSC 92-2416-H-004-004, NSC 93-2416-H-004-015, and NSC 43028F.

$\{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2), j = 1, \dots, p\}$.

$$w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) = \mathbf{T}. \quad (3)$$

Hereafter, ${}_a\mathbf{R}_b \equiv \mathbf{a} - \frac{\mathbf{b}^T \mathbf{a}}{\|\mathbf{b}\|} \frac{\mathbf{b}}{\|\mathbf{b}\|}$ denotes the residual of vector \mathbf{a} regarding \mathbf{b} after the part parallel

with the vector \mathbf{b} has been taken away. Furthermore, ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} \equiv {}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}$

$-\frac{\mathbf{a}^T {}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|} \frac{{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}}{\|{}_{\mathbf{b}^k} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k-1}\}}\|}$ denotes the residual of vector \mathbf{a} regarding the ordered sequence of

vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$ after the parts parallel with vectors \mathbf{b}^j 's have been taken away.

Note that ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} = \mathbf{a} - \frac{\mathbf{a}^T \mathbf{b}^1}{\|\mathbf{b}^1\|} \frac{\mathbf{b}^1}{\|\mathbf{b}^1\|} - \sum_{j=2}^k \frac{\mathbf{a}^T {}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} \frac{{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|}$ and thus $\mathbf{a} = \frac{\mathbf{a}^T \mathbf{b}^1}{\|\mathbf{b}^1\|}$

$\frac{\mathbf{b}^1}{\|\mathbf{b}^1\|} + \sum_{j=2}^k \frac{\mathbf{a}^T {}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} \frac{{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}}{\|{}_{\mathbf{b}^j} \mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{j-1}\}}\|} + {}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$. Furthermore, we have **Lemma 1** below, which

lists some properties associated with the above proposed residual vector.

Lemma 1: (i) If \mathbf{a} is linearly dependent with \mathbf{b} , then ${}_a\mathbf{R}_b = \mathbf{0}$. Similarly, if \mathbf{a} can be linearly represented by the set of vectors $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$,¹ then ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} = \mathbf{0}$. (ii) ${}_a\mathbf{R}_b$ is orthogonal to \mathbf{b} . Similarly, ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}}$ is orthogonal to all vectors in the set $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$. (iii) If ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^k\}} = \mathbf{0}$, then ${}_a\mathbf{R}_{\{\mathbf{b}^1, \dots, \mathbf{b}^{k+1}\}} = \mathbf{0}$.

3. The Proposed Constructive Procedure

Table 1 presents the proposed procedure for constructing a SLFN appropriate for fitting the mapping embedded in $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$.

Table 1. The proposed deterministic procedure for constructing an appropriate SLFN for the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$. $\mathbf{T} \equiv (t^1, t^2, \dots, t^N)^T$ and $\mathbf{v}(w_0, \mathbf{w}) \equiv (a^1(w_0, \mathbf{w}), a^2(w_0, \mathbf{w}), \dots, a^N(w_0, \mathbf{w}))^T$ is a $N \times 1$ vector with the c^{th} component being $a^c(w_0, \mathbf{w}) \equiv \tanh(\mathbf{w}^T \mathbf{x}^c + w_0)$, in which $\mathbf{w} \equiv (w_1, w_2, \dots, w_m)^T$.

Step 1: Calculate $\mathbf{u}^0 = \frac{\mathbf{1}}{\|\mathbf{1}\|}$ and ${}_T\mathbf{R}_{\mathbf{u}^0} \equiv \mathbf{T} - (\mathbf{u}^0)^T \mathbf{T} \mathbf{u}^0$.

Step 2: If ${}_T\mathbf{R}_{\mathbf{u}^0} = \mathbf{0}$, then (i) claim that the fitting job requests no hidden node; (ii) set the bias of output

node as $\frac{\mathbf{T}^T \mathbf{u}^0}{\mathbf{1}^T \mathbf{u}^0}$ and the weight vector between the output node and the input layer as $\mathbf{0}$; and (iii)

stop.

Step 3: If ${}_T\mathbf{R}_{\mathbf{u}^0} \neq \mathbf{0}$, then solve $\min_{w_0, \mathbf{w}} \|{}_T\mathbf{R}_{\mathbf{v}(w_0, \mathbf{w})}\|^2$ and let w_0^* and \mathbf{w}^* be the obtained optimal solution.

¹ Namely, \mathbf{a} is in the (sub-)space spanned by $\{\mathbf{b}^1, \dots, \mathbf{b}^k\}$.

Then set $w_{10}^2 = w_0^*$, $\mathbf{w}_1^2 = \mathbf{w}^*$, $\mathbf{u}^1 \equiv \frac{\mathbf{v}(w_{10}^2, \mathbf{w}_1^2)}{\|\mathbf{v}(w_{10}^2, \mathbf{w}_1^2)\|}$, and $p = 1$.

Step 4: Calculate $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}$.

Step 5: If $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{0}$, then (i) claim that the fitting job requests p hidden nodes with the bias and weights of the j^{th} hidden node being the above w_{j0}^2 and \mathbf{w}_j^2 , $\forall j = 1, \dots, p$; (ii) set $w_0^3 \equiv 0$,

$$w_p^3 \equiv \frac{\mathbf{T}^T \mathbf{u}^p}{\mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)^T \mathbf{u}^p}, \text{ and } w_j^3 \equiv \frac{(\mathbf{T} - \sum_{l=j+1}^p w_l^3 \mathbf{v}(w_{l0}^2, \mathbf{w}_l^2))^T \mathbf{u}^j}{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)^T \mathbf{u}^j} \text{ for all } 1 \leq j \leq p-1; \text{ and}$$

(iii) stop.

Step 6: If $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}} \neq \mathbf{0}$, then calculate $\mathbf{u}^0 = \frac{\mathbf{1} \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{1} \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$ and $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} \equiv \tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}} - (\mathbf{u}^0)^T \mathbf{T} \mathbf{u}^0$.

Step 7: If $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} = \mathbf{0}$, then (i) claim that the fitting job requests p hidden nodes with the bias and weights of the j^{th} hidden node being the above w_{j0}^2 and \mathbf{w}_j^2 , $\forall j = 1, \dots, p$; (ii) set

$$w_0^3 \equiv \frac{\mathbf{T}^T \mathbf{u}^0}{\mathbf{1}^T \mathbf{u}^0}, w_p^3 \equiv \frac{(\mathbf{T} - w_0^3 \mathbf{1})^T \mathbf{u}^p}{\mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)^T \mathbf{u}^p}, \text{ and } w_j^3 \equiv \frac{(\mathbf{T} - w_0^3 \mathbf{1} - \sum_{l=j+1}^p w_l^3 \mathbf{v}(w_{l0}^2, \mathbf{w}_l^2))^T \mathbf{u}^j}{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)^T \mathbf{u}^j}$$

for all $1 \leq j \leq p-1$; and (iii) stop.

Step 8: If $\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} \neq \mathbf{0}$, then solve $\min_{w_0, \mathbf{w}} \|\tau\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{v}(w_0, \mathbf{w})\}}\|^2$ and let w_0^* and \mathbf{w}^* be the obtained optimal

solution. Then (i) set $w_{p+1,0}^2 = w_0^*$, $\mathbf{w}_{p+1}^2 = \mathbf{w}^*$, $\mathbf{u}^{p+1} \equiv \frac{\mathbf{v}(w_{p+1,0}^2, \mathbf{w}_{p+1}^2) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{v}(w_{p+1,0}^2, \mathbf{w}_{p+1}^2) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$, set $p+1 \rightarrow p$; and

(ii) go to Step 4.

4. The Correctness of the Proposed Procedure

We now prove that the procedure stated in Table 1 creates an appropriate SLFN that meets the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$ without error.

Lemma 2 below is obvious from system (3) and the definition of the residual vector. **Lemma 2** states the necessary condition of a SLFN with p hidden nodes appropriate for the mapping requirement of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$. Note that, from Lemma 1 (iii), $\tau\mathbf{R}_{\{\mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} = \mathbf{0}$ results in $\tau\mathbf{R}_{\{\mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{1}\}} = \mathbf{0}$. Thus the condition stated in Steps 2,

5, and 7 are suitable stopping criteria.

Lemma 2: Regarding the mapping samples of $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$, the SLFN with p hidden nodes is appropriate for the mapping requirement if the associated $\mathbf{T}\mathbf{R}_{\{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{1}\}}$ equals $\mathbf{0}$.

Namely, if $\mathbf{T}\mathbf{R}_{\mathbf{u}^0} = \mathbf{0}$ at Step 2 of Table 1, in which $\mathbf{u}^0 = \frac{\mathbf{1}}{\|\mathbf{1}\|}$, then $\mathbf{T} = \frac{\mathbf{T}^T \mathbf{u}^0}{\mathbf{1}^T \mathbf{u}^0}$ and the mapping requirement asks for no hidden node. Furthermore, because $\mathbf{u}^1 \equiv \frac{\mathbf{v}(w_{10}^2, \mathbf{w}_1^2)}{\|\mathbf{v}(w_{10}^2, \mathbf{w}_1^2)\|}$, $\mathbf{u}^j \equiv \frac{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^{j-1}\}}}{\|\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^{j-1}\}}\|}$ for all $j = 2, \dots, p$, and $\mathbf{u}^0 \equiv \frac{\mathbf{1} \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{1} \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$, $\mathbf{T}\mathbf{R}_{\{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)\}} = \mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}$ and $\mathbf{T}\mathbf{R}_{\{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2), \dots, \mathbf{v}(w_{p0}^2, \mathbf{w}_p^2), \mathbf{1}\}} = \mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}}$. Thus, if $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}} = \mathbf{0}$ at Step 5 of Table 1, then $\mathbf{T} = \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)$ and the mapping requirement asks for p hidden nodes. If $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} = \mathbf{0}$ at Step 7 of Table 1, then $\mathbf{T} = w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)$ and the mapping requirement asks for p hidden nodes.

The following **Lemmas 3** and **4** show that the proposed procedure will generate a sequence of orthonormal vectors.

Lemma 3: If $p < N-1$, then there always exist w_0 and \mathbf{w} such that $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} \neq \mathbf{0}$.

Proof of Lemma 3: Suppose $p < N-1$ and there are no w_0 and \mathbf{w} such that $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} \neq \mathbf{0}$. In other words, $p < N-1$ and there are no w_0 and \mathbf{w} such that $\mathbf{v}(w_0, \mathbf{w})$ is linearly independent with $\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}$. This contradicts with the statement of [4] that, for any given $\{\mathbf{x}^c \mid c = 1, \dots, N\}$, there exists a set of $\{(w_{j0}^2, \mathbf{w}_j^2), j = 1, \dots, N-1\}$ such that the associated square matrix $(\mathbf{1}, \mathbf{v}(w_{10}^2, \mathbf{w}_1^2), \dots, \mathbf{v}(w_{N-1,0}^2, \mathbf{w}_{N-1}^2))$ is invertible. **Q.E.D.**

Lemma 4: When the proposed procedure stops at some p , $1 \leq p \leq N-1$, all vectors in the set of $\{\mathbf{u}^1, \dots, \mathbf{u}^p\}$ or $\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}$ are orthonormal.

Proof of Lemma 4: It is trivial for the case of $p = 1$ and $\{\mathbf{u}^1\}$. As for the case of $p = 1$ and $\{\mathbf{u}^1, \mathbf{u}^0\}$, from steps 6 and 7 of Table 1, $\mathbf{T}\mathbf{R}_{\mathbf{u}^1} \neq \mathbf{0}$, $\mathbf{u}^0 \equiv \frac{\mathbf{1} \mathbf{R}_{\mathbf{u}^1}}{\|\mathbf{1} \mathbf{R}_{\mathbf{u}^1}\|}$, and $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \mathbf{u}^0\}} = \mathbf{0}$. Thus \mathbf{u}^0 is non-zero and, from

Lemma 1 (ii), orthogonal to \mathbf{u}^1 . Namely, $\{\mathbf{u}^1, \mathbf{u}^0\}$ is an orthonormal set.

Now consider the case of any p with $2 \leq p \leq N-1$. Then, for each $1 \leq j \leq p-1$, from **Lemma 3**, there exists w_0 and \mathbf{w} such that $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{u}^0\}}$ is not zero and thus, from **Lemma 1** (iii), $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}$ is not zero. Furthermore, as stated in Step 8 of Table 1, $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{u}^0\}} \neq \mathbf{0}$. Thus $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}} \neq \mathbf{0}$, there exists a non-zero $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}$ that $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}} = \mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{v}(w_0, \mathbf{w})\}}$ + $\frac{\mathbf{T}^T \mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}}{\|\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}\|^2} \mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}$, and, from **Lemma 1** (ii), $\mathbf{v}(w_0, \mathbf{w}) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}$ and $\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{v}(w_0, \mathbf{w})\}}$ are orthogonal to each other. Namely, $\|\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}\|^2 \geq \|\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{v}(w_0, \mathbf{w})\}}\|^2$. Therefore, the optimization of $\min_{w_0, \mathbf{w}} \|\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{v}(w_0, \mathbf{w})\}}\|^2$ leads to a non-zero $\mathbf{v}(w_0^*, \mathbf{w}^*) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}$, in which $(w_0^*, \mathbf{w}^*) \equiv \arg(\min_{w_0, \mathbf{w}} \|\mathbf{T}\mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j, \mathbf{v}(w_0, \mathbf{w})\}}\|^2)$. Thus $\mathbf{u}^{j+1} \equiv \frac{\mathbf{v}(w_0^*, \mathbf{w}^*) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}}{\|\mathbf{v}(w_0^*, \mathbf{w}^*) \mathbf{R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^j\}}\|}$ is a unit vector and, from **Lemma 1** (ii), orthogonal to the set of vectors $\{\mathbf{u}^1, \dots, \mathbf{u}^j\}$. So, all vectors in the generated set $\{\mathbf{u}^1, \dots, \mathbf{u}^p\}$ at Step 8 are orthonormal.

From Steps 6 and 7 of Table 1, if $\mathbf{TR}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}} \neq \mathbf{0}$, $\mathbf{u}^0 = \frac{\mathbf{1R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}}{\|\mathbf{1R}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p\}}\|}$ and $\mathbf{TR}_{\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}} = 0$, then \mathbf{u}^0 is non-zero and,

from **Lemma 1** (ii), orthogonal to all vectors in the set $\{\mathbf{u}^1, \dots, \mathbf{u}^p\}$. Namely, $\{\mathbf{u}^1, \dots, \mathbf{u}^p, \mathbf{u}^0\}$ is an orthonormal set.

Q.E.D.

If $\mathbf{T} = \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)$, then $\mathbf{T}^T \mathbf{u}^k = (\sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2))^T \mathbf{u}^k = \sum_{j=k}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)^T \mathbf{u}^k$, since $\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)$ is in the

subspace spanned by $\{\mathbf{u}^1, \dots, \mathbf{u}^j\}$ for all $j = 1, \dots, p$ and $(\mathbf{u}^i)^T \mathbf{u}^j \neq 0 \forall i \neq j$. Thus, as stated in Step 5 of Table 1,

$$w_p^3 \equiv \frac{\mathbf{T}^T \mathbf{u}^p}{\mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)^T \mathbf{u}^p}, \quad w_j^3 \equiv \frac{(\mathbf{T} - \sum_{l=j+1}^p w_l^3 \mathbf{v}(w_{l0}^2, \mathbf{w}_l^2))^T \mathbf{u}^j}{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)^T \mathbf{u}^j} \quad \text{for all } 1 \leq j \leq p-1, \text{ and } w_0^3 \equiv 0. \text{ Similarly, if } \mathbf{T} = w_0^3 \mathbf{1}$$

+ $\sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)$, then $\mathbf{T}^T \mathbf{u}^k = (w_0^3 \mathbf{1} + \sum_{j=1}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2))^T \mathbf{u}^k = (w_0^3 \mathbf{1} + \sum_{j=k}^p w_j^3 \mathbf{v}(w_{j0}^2, \mathbf{w}_j^2))^T \mathbf{u}^k$. Therefore, as

$$\text{stated in Step 7 of Table 1, } w_0^3 \equiv \frac{\mathbf{T}^T \mathbf{u}^0}{\mathbf{1}^T \mathbf{u}^0}, \quad w_p^3 \equiv \frac{(\mathbf{T} - w_0^3 \mathbf{1})^T \mathbf{u}^p}{\mathbf{v}(w_{p0}^2, \mathbf{w}_p^2)^T \mathbf{u}^p}, \text{ and } w_j^3 \equiv \frac{(\mathbf{T} - w_0^3 \mathbf{1} - \sum_{l=j+1}^p w_l^3 \mathbf{v}(w_{l0}^2, \mathbf{w}_l^2))^T \mathbf{u}^j}{\mathbf{v}(w_{j0}^2, \mathbf{w}_j^2)^T \mathbf{u}^j} \quad \text{for all } 1 \leq$$

$j \leq p-1$.

References

1. M. Me'zard and J. Nadal, *Learning in feedforward layered networks: The tiling algorithm*. *Journal of Physics*. **A 22**, 2191– 2204(1989).
2. S. Fahlman and C. Lebiere, *The Cascade-Correlation Learning Architecture*. In Touretzky, D. (Eds.), *Advances in Neural Information Processing Systems II (Denver, 1989)* Morgan Kaufmann, San Mateo (1990).
3. M. Frean, *The Upstart Algorithm: A Method for Constructing and Training Feedforward Neural Networks*. *Neural Computation*. **2**, 198– 209(1990).
4. G. Huang and H. Babri, *Upper bounds on the number of hidden neurons in feedforward networks with arbitrary bounded nonlinear activation functions*. *IEEE Transactions on Neural Networks*. **9**, 224– 229(1998).