



Cholesky factorization based online regularized and kernelized extreme learning machines with forgetting mechanism



Xin-Ran Zhou, Chun-Sheng Wang*

School of Information Science and Engineering, Central South University, Changsha, China

ARTICLE INFO

Article history:

Received 25 June 2015

Received in revised form

21 September 2015

Accepted 11 October 2015

Communicated by G.-B. Huang

Available online 23 October 2015

Keywords:

Extreme learning machine

Online sequential learning algorithms

Forgetting mechanism

Cholesky decomposition

ABSTRACT

In this paper, we propose two alternative schemes of fast online sequential extreme learning machine (ELM) for training the single hidden-layer feedforward neural networks (SLFN), termed as Cholesky factorization based online regularized ELM with forgetting mechanism (CF-FORELM) and Cholesky factorization based online kernelized ELM with forgetting mechanism (CF-FOKELM). First, the solutions of regularized ELM (RELM) and kernelized ELM (KELM) using the matrix Cholesky factorization are introduced; then the recursive method for calculating Cholesky factor of involved matrix in RELM and KELM is designed when RELM and KELM are applied to train SLFN online; consequently, the CF-FORELM and CF-FOKELM are obtained. The numerical simulation results show CF-FORELM demands less computational burden than Dynamic Regression ELM (DR-ELM), and CF-FOKELM also owns higher computational efficiency than both FOKELM and online sequential ELM with kernels (OS-ELMK), and CF-FORELM is less sensitive to model parameters than CF-FOKELM.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Lots of research work has shown that the single hidden-layer feedforward neural networks (SLFN) can approximate any function and form decision boundaries with arbitrary shapes if the activation function is chosen properly [1–3]. To train SLFN fast, Huang et al. proposed a learning algorithm called extreme learning machine (ELM), which randomly assigns the hidden nodes parameters (the input weights and hidden layer biases of additive networks or the centers and impact factors of RBF networks) and then determines the output weights by the Moore–Penrose generalized inverse [4,5]. The original ELM is batch learning algorithm. Over the past decade, owing to its high efficiency, ELM has gained increasing interest from various research fields and its variants have been designed for different purposes under general or specific conditions [6–10].

In order to train SLFN to achieve better generalization ability, Regularized ELM (RELM) [11–13], which is equivalent to the constrained optimization based ELM [14,15] mathematically, was presented by researchers. Besides, calculating inverse of invertible matrix in RELM is less complex than calculating Moore–Penrose generalized inverse in ELM; the RELM can greatly reduce the randomness effect in ELM [16].

In order to train SLFN to capture late characteristics of identified systems, especially nonstationary time series or time-varying systems where training data have timeliness, Zhang and Wang, using the block matrix inverse formula [17], proposed Dynamic Regression ELM (DR-ELM) [18], that is, Local ELM (LELM) [19]. Recently, Zhou et al. presented online RELM with forgetting mechanism (FORELM) [20] invoking Sherman–Morrison formula [17]. Both DR-ELM and FORELM study a given number of samples, i.e., those samples during the sliding time window.

If the feature mapping in SLFN is unknown, the kernelized ELM (KELM) can be constructed [14,15]. In order to train SLFN to identify time-varying or nonstationary systems, Zhou et al. presented online kernelized ELM with forgetting mechanism (FOKELM) [20], and Wang and Han proposed an online sequential extreme learning machine with kernels (OS-ELMK) [21]. Both FOKELM and OS-ELMK also study samples during the sliding time window only. We do not need to select the number of hidden nodes for these KELMs.

The time efficiency of learning algorithm is an important index in general. In realtime applications, such as stock forecast, modelling for controlled objects, signal processing, the computational efficiency of online training algorithm for SLFN is a crucial factor. Hence, in order to accelerate DR-ELM, FOKELM and OS-ELMK, we design Cholesky factorization based online regularized ELM with forgetting mechanism (CF-FORELM) and Cholesky factorization based online kernelized ELM with forgetting mechanism (CF-FOKELM). The numerical experiments show CF-FORELM runs faster than DR-ELM, and CF-FOKELM also owns higher computational

* Corresponding author.

E-mail address: wangcsh2015@126.com (C.-S. Wang).

efficiency than both FOKELM and OS-ELMK. Furthermore, we compare parameters sensitivity between CF-FORELM with the sigmoidal additive function and CF-FOKELM with Gaussian kernel function, contrast their speed, and offer suggestions to choose CF-FORELM or CF-FOKELM.

It should be noticed that herein CF-FORELM, like FORELM, tunes the output weights of SLFN due to addition and deletion of the samples one by one, viz., learns and forgets samples sequentially, and network architecture is fixed. It is completely different from those offline incremental RELM [22–24] seeking optimal network architecture by adding hidden nodes one by one, and learning the data in batch mode.

The rest of this paper is organized as follows. Section 2 gives a brief review of the basic concepts and related work of RELM and KELM. Section 3 introduces solutions of RELM and KELM by using the Cholesky factorization. Section 4 proposes two new online learning algorithms, viz., CF-FORELM and CF-FOKELM. Performance evaluation is conducted in Section 5. Finally, conclusions and discussion are given in Section 6.

2. Brief review of the RELM and KELM

For simplicity, ELM based learning algorithm for SLFN with multiple input single output is discussed.

The output of a SLFN with L hidden nodes (additive or RBF nodes) can be represented by

$$f(\mathbf{x}) = \sum_{i=1}^L \beta_i G(\mathbf{a}_i, b_i, \mathbf{x}) = \mathbf{h}(\mathbf{x})\boldsymbol{\beta} \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{a}_i \in \mathbb{R}^n, \quad (1)$$

where \mathbf{a}_i and b_i are the learning parameters of hidden nodes, and $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_L]^T$ is the vector of the output weights, and $G(\mathbf{a}_i, b_i, \mathbf{x})$ denotes the output of the i -th hidden node with respect to the input \mathbf{x} , i.e., activation function. $\mathbf{h}(\mathbf{x}) = [G(\mathbf{a}_1, b_1, \mathbf{x}), G(\mathbf{a}_2, b_2, \mathbf{x}), \dots, G(\mathbf{a}_L, b_L, \mathbf{x})]$ is a feature mapping from the n -dimensional input space to the L -dimensional hidden-layer feature space. In ELM, \mathbf{a}_i and b_i are randomly determined firstly.

For a given set of distinct training data $\{(\mathbf{x}_i, y_i)\}_i^N \subset \mathbb{R}^n \times \mathbb{R}$, where \mathbf{x}_i is an n dimension input vector and y_i is the corresponding scalar observation, The RELM, i.e., constrained optimization based ELM, can be formulated as

$$\begin{aligned} \text{Minimize : } & L_{P_{ELM}} = \frac{1}{2} \|\boldsymbol{\beta}\|^2 + c \frac{1}{2} \|\boldsymbol{\xi}\|^2 \\ \text{Subject to : } & \mathbf{H}\boldsymbol{\beta} = \mathbf{Y} - \boldsymbol{\xi} \end{aligned} \quad (2)$$

where $\boldsymbol{\xi} = [\xi_1, \xi_2, \dots, \xi_N]^T$ denotes the training error. $\mathbf{Y} = [y_1, y_2, \dots, y_N]^T$ indicates the target value of all the samples. $\mathbf{H}_{N \times L} = [\mathbf{h}(\mathbf{x}_1)^T, \mathbf{h}(\mathbf{x}_2)^T, \dots, \mathbf{h}(\mathbf{x}_N)^T]^T$ is the mapping matrix for the inputs of all the samples. c is the regularization parameter (a positive constant).

Based on the KKT theorem, the constrained optimization of (2) can be transferred to the following dual optimization problem:

$$L_{D_{ELM}} = \frac{1}{2} \|\boldsymbol{\beta}\|^2 + c \frac{1}{2} \|\boldsymbol{\xi}\|^2 - \boldsymbol{\alpha}^T (\mathbf{H}\boldsymbol{\beta} - \mathbf{Y} + \boldsymbol{\xi}) \quad (3)$$

where $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N]^T$ is the Lagrange multipliers vector. Using KKT optimality conditions, the following equations can be gotten:

$$\frac{\partial L_{D_{ELM}}}{\partial \boldsymbol{\beta}} = \mathbf{0} \rightarrow \boldsymbol{\beta} = \mathbf{H}^T \boldsymbol{\alpha} \quad (4)$$

$$\frac{\partial L_{D_{ELM}}}{\partial \boldsymbol{\xi}} = \mathbf{0} \rightarrow \boldsymbol{\alpha} = c \boldsymbol{\xi} \quad (5)$$

$$\frac{\partial L_{D_{ELM}}}{\partial \boldsymbol{\alpha}} = \mathbf{0} \rightarrow \mathbf{H}\boldsymbol{\beta} - \mathbf{Y} + \boldsymbol{\xi} = \mathbf{0} \quad (6)$$

Ultimately, $\boldsymbol{\beta}$ can be obtained as follows [11,14,15]:

$$\boldsymbol{\beta} = (c^{-1} \mathbf{I} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y} \quad (7)$$

or

$$\boldsymbol{\beta} = \mathbf{H}^T (c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T)^{-1} \mathbf{Y} \quad (8)$$

In order to reduce computational costs, when $N > L$, one may prefer to apply solution (7), and when $N < L$, one may prefer to apply solution (8).

Invoking Eq. (8), the output of SLFN can be obtained as

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{h}(\mathbf{x})\boldsymbol{\beta} = \mathbf{h}(\mathbf{x})\mathbf{H}^T (c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T)^{-1} \mathbf{Y} = \mathbf{h}(\mathbf{x})[\mathbf{h}(\mathbf{x}_1)^T, \mathbf{h}(\mathbf{x}_2)^T, \dots, \mathbf{h}(\mathbf{x}_N)^T] \\ &\quad (c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T)^{-1} \mathbf{Y} \\ &= [\mathbf{h}(\mathbf{x})\mathbf{h}(\mathbf{x}_1)^T, \mathbf{h}(\mathbf{x})\mathbf{h}(\mathbf{x}_2)^T, \dots, \mathbf{h}(\mathbf{x})\mathbf{h}(\mathbf{x}_N)^T] \mathbf{g} \\ &= \sum_{i=1}^N g_i \mathbf{h}(\mathbf{x})\mathbf{h}(\mathbf{x}_i)^T \end{aligned} \quad (9)$$

where $\mathbf{g} = [g_1, g_2, \dots, g_N]^T$. Obviously, \mathbf{g} can be calculated by solving the following equation:

$$(c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T) \mathbf{g} = \mathbf{Y} \quad (10)$$

If the feature mapping $\mathbf{h}(\mathbf{x})$ is unknown, one can apply Mercer's conditions on RELM. The kernel matrix is defined as $\boldsymbol{\Omega} = \mathbf{H}\mathbf{H}^T$ and $\Omega_{ij} = \mathbf{h}(\mathbf{x}_i)\mathbf{h}(\mathbf{x}_j)^T = K(\mathbf{x}_i, \mathbf{x}_j)$. Then, the output of SLFN by kernel based RELM can be given as

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{h}(\mathbf{x})\boldsymbol{\beta} = \mathbf{h}(\mathbf{x})\mathbf{H}^T (c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T)^{-1} \mathbf{Y} \\ &= [K(\mathbf{x}, \mathbf{x}_1), K(\mathbf{x}, \mathbf{x}_2), \dots, K(\mathbf{x}, \mathbf{x}_N)] (c^{-1} \mathbf{I} + \boldsymbol{\Omega})^{-1} \mathbf{Y} \\ &= [K(\mathbf{x}, \mathbf{x}_1), K(\mathbf{x}, \mathbf{x}_2), \dots, K(\mathbf{x}, \mathbf{x}_N)] \mathbf{g} \\ &= \sum_{i=1}^N g_i K(\mathbf{x}, \mathbf{x}_i) \end{aligned} \quad (11)$$

here \mathbf{g} can be calculated by solving the following equation:

$$(c^{-1} \mathbf{I} + \boldsymbol{\Omega}) \mathbf{g} = \mathbf{Y} \quad (12)$$

3. The solutions of RELM and KELM by the Cholesky factorization

Theorem 1. The matrix $c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T$ is a symmetrical positive definition matrix.

Proof. Symmetry. Apparently,

$$(c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T)^T = c^{-1} \mathbf{I} + (\mathbf{H}\mathbf{H}^T)^T = c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T \quad (13)$$

Positive definition property. For any $\boldsymbol{\zeta} = [\zeta_1, \zeta_2, \dots, \zeta_N] \neq \mathbf{0}$, it holds that

$$\boldsymbol{\zeta} c^{-1} \mathbf{I} \boldsymbol{\zeta}^T = c^{-1} \sum_{i=1}^N \zeta_i^2 > 0. \quad (14)$$

□

Additionally,

$$\begin{aligned} \boldsymbol{\zeta} \mathbf{H}\mathbf{H}^T \boldsymbol{\zeta}^T &= \boldsymbol{\zeta} \mathbf{H} (\boldsymbol{\zeta} \mathbf{H})^T = \left(\sum_{i=1}^N \zeta_i G(\mathbf{a}_1, b_1, \mathbf{x}_i) \right)^2 + \left(\sum_{i=1}^N \zeta_i G(\mathbf{a}_2, b_2, \mathbf{x}_i) \right)^2 \\ &\quad + \dots + \left(\sum_{i=1}^N \zeta_i G(\mathbf{a}_L, b_L, \mathbf{x}_i) \right)^2 \geq 0 \end{aligned} \quad (15)$$

Now that $c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T$ is symmetric positive definite matrix. Denote $\mathbf{B} = c^{-1} \mathbf{I} + \mathbf{H}\mathbf{H}^T$, then, \mathbf{B} can be factorized in Cholesky form, i.e., $\mathbf{B} = \mathbf{U}^T \mathbf{U}$, where \mathbf{U} is an upper triangular matrix. Then \mathbf{U} can be calculated as following formulas:

$$u_{ii} = \left(b_{ii} - \sum_{d=1}^{i-1} u_{di}^2 \right)^{1/2}, \quad i = 1, \dots, N \quad (16)$$

Download English Version:

<https://daneshyari.com/en/article/411648>

Download Persian Version:

<https://daneshyari.com/article/411648>

[Daneshyari.com](https://daneshyari.com)