



Extreme spectral regression for efficient regularized subspace learning



Bing Liu, Shi-Xiong Xia*, Fan-Rong Meng, Yong Zhou

School of Computer Science and Technology, China University of Mining, and Technology, Jiangsu, Xuzhou 221116, China

ARTICLE INFO

Article history:

Received 6 August 2013
 Received in revised form
 26 September 2013
 Accepted 30 September 2013
 Available online 30 September 2014

Keywords:

Spectral regression
 Extreme learning machine (ELM)
 Dimensionality reduction
 Out-of-sample extension

ABSTRACT

Traditional manifold learning algorithms, such as Locally Linear Embedding, Isomap and Laplacian Eigenmap, only provide the embedding results of training samples. Although many extensions of these approaches try to solve the out-of-sample extension problem, their computations cannot avoid eigen-decomposition of dense matrices which is expensive in both time and memory. To solve this problem, spectral regression (SR) casts the problem of learning an embedding function into a regression framework. Motivated by the effectiveness of extreme learning machine (ELM), in this paper, we solve the out-of-sample extension problem by seeking an embedding function in ELM feature space. An extreme spectral regression (ESR) algorithm is proposed to speed up kernel-based SR (KSR) further. In addition, it is proved that ESR is an approximation of KSR. Similar to SR, the proposed ESR algorithm can be performed in supervised, unsupervised and semi-supervised situation. Experimental results on classification and semi-supervised classification demonstrate the effectiveness and efficiency of our algorithm.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

Dimensionality reduction has been a key problem in many fields of information processing, such as machine learning, data mining, information retrieval, and pattern recognition. Practical algorithms usually behave badly when faced with many unnecessary features. A common way to attempt to resolve this problem is to use dimensionality reduction (DR) techniques, which include unsupervised, supervised and semi-supervised DR due to different assumptions about the data distribution or the availability of the data labeling.

Principal component analysis (PCA) [1] is one of the most popular unsupervised DR techniques, which finds a linear mapping by maximizing the projected variances. If the data is embedded in a linear subspace, PCA is guaranteed to discover the dimensionality of the subspace and produces a compact representation. In order to handle the data sampled from a nonlinear low dimensional manifold, many manifold learning techniques, such as ISOMAP [2], Locally Linear Embedding (LLE) [3] and Laplacian Eigenmap [4] have been proposed which reduce the dimensionality of a fixed training set in a way that can maximally preserve certain inter-point relationships. Given the data of each class have a Gaussian distribution, Linear discriminant analysis (LDA) first constructs the between-class scatter and the within-class scatter matrices by virtue of labeled data, then simultaneously maximizes the between-class scatter and minimizes the

within-class scatter to obtain a projection. Alternatively, marginal Fisher analysis (MFA) [5] and local discriminant embedding (LDE) [6] exploit the assumption that the data of each class spread as a submanifold, and seek a discriminant embedding over these submanifolds. Semi-supervised DR, such as semi-supervised discriminant analysis (SDA) [7], utilizes partially labeled data while preserving the intrinsic geometric structures of the remaining. Generally, linear DR methods mentioned above can be kernelized into nonlinear ones. As shown in [8–10], the kernelized versions can achieve significant improvements.

Although some modified methods explicitly require an embedding function either linear or in the Reproducing Kernel Hilbert Space (RKHS) when minimizing the objective function [11,12], the computation of these methods involves eigen-decomposition of dense matrices which is expensive in both time and memory. It is almost infeasible to apply these approaches on large data sets. Spectral regression (SR), which is fundamentally based on regression and spectral graph analysis [13–16], casts the problem of learning an embedding function into a regression framework. In this method, an affinity graph over both labeled and unlabeled points is first constructed to discover the intrinsic discriminant structure in the data. The responses for both labeled and unlabeled points are then obtained by means of this graph. Finally, the ordinary regression is applied for learning the embedding function. Thus, SR avoids eigen-decomposition of dense matrices. Moreover, it can be performed either in supervised, unsupervised or semi-supervised situation.

Kernel SR (KSR) is the kernelized version of SR in the reproducing kernel Hilbert space (RKHS) into which data points are

* Corresponding author.

E-mail address: xiasx@cumt.edu.cn (S.-X. Xia).

Table 1
Notations.

Notations	Descriptions
\mathbb{R}^d	the input d -dimensional Euclidean space
n	the number of total training data points
m	the number of classes that the samples belong to
\mathbf{X}	$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ is the training data matrix
$k(\mathbf{x}, \mathbf{y})$	Kernel function of variables \mathbf{x} and \mathbf{y}
\mathbf{K}	Kernel matrix $\mathbf{K} = \{k(\mathbf{x}_i, \mathbf{x}_j)\} \in \mathbb{R}^{n \times n}$
$\ \cdot\ $	norm in the Hilbert space \mathcal{H}
\mathbf{H}	the hidden-layer output matrix
β	the vector of the output weights between the hidden layer of L nodes and the output node
L	the graph Laplacian matrix

mapped. For large data sets, the computation of KSR involves an inversion of the kernel matrix of the training data, which is time-consuming in applications. To solve this problem, motivated by the fast learning speed of extreme learning machine (ELM), we solve the out-of-sample extension problem by seeking an embedding function in ELM feature space instead of RKHS. Thus, the final regression problem can be solved effectively by the ELM algorithm, which uses L_2 regularization to solve the over-fitting problem. The regularization methods for solving the regression problem of ELM have been studied extensively. Miche et al. in [17,18] proposed a double regularized ELM algorithm, which uses a cascade of two regularization penalties: first a L_1 penalty to rank the neurons of the hidden layer, followed by a L_2 penalty to prune the network accordingly. In this paper, we propose an extreme spectral regression (ESR) algorithm, which is based on ELM only using L_2 regularization for simplicity. In this case, it is proved that ESR is an approximation of KSR by discretizing the embedding functions in RKHS. Experimental results on classification and semi-supervised classification demonstrate the effectiveness and efficiency of our algorithm.

The paper is structured as follows. In Section 2, we briefly introduce the extreme learning machine model. Our extreme spectral regression algorithm is introduced in Section 3. In Section 4, we provide a theoretical and computational complexity analysis of our algorithm respectively. The experimental results are presented in Section 5. Finally, we give the related conclusions in Section 6. In order to avoid confusion, we give a list of the main notations used in this paper in Table 1.

2. Extreme learning machine

The output function of ELM for generalized SLFNs in the case of one output node is

$$f_L(\mathbf{x}) = \sum_{i=1}^L \beta_i h_i(\mathbf{x}) = \mathbf{h}(\mathbf{x})\beta, \quad (1)$$

where $\beta = [\beta_1, \dots, \beta_L]^T$ is the vector of the output weights between the hidden layer of L nodes and the output node, and $\mathbf{h}(\mathbf{x}) = [h_1(\mathbf{x}), \dots, h_L(\mathbf{x})]$ is the output (row) vector of the hidden layer with respect to the input \mathbf{x} . In fact, $\mathbf{h}(\mathbf{x})$ maps the data from the d -dimensional input space to the L -dimensional hidden-layer feature space (ELM feature space) \mathbf{H} . ELM is to minimize the training error as well as the norm of the output weights [18–20]

$$\min_{\beta} \frac{C}{2} \|\mathbf{H}\beta - \mathbf{T}\|^2 + \frac{1}{2} \|\beta\|^2, \quad (2)$$

where C is a tradeoff parameter between the complexity and fitness of the decision function and \mathbf{H} is the hidden-layer output

matrix denoted by

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}(\mathbf{x}_1) \\ \mathbf{h}(\mathbf{x}_2) \\ \vdots \\ \mathbf{h}(\mathbf{x}_n) \end{bmatrix} = \begin{bmatrix} h_1(\mathbf{x}_1) & \dots & h_L(\mathbf{x}_1) \\ h_1(\mathbf{x}_2) & \dots & h_L(\mathbf{x}_2) \\ \vdots & \vdots & \vdots \\ h_1(\mathbf{x}_n) & \dots & h_L(\mathbf{x}_n) \end{bmatrix}. \quad (3)$$

Similar to support vector machine (SVM), to minimize the norm of the output weights $\|\beta\|$ is actually to maximize the distance of the separating margins of the two different classes in the ELM feature space: $2/\|\beta\|$, which actually controls the complexity of the function in the ELM feature space.

For completeness, we briefly introduce the multiclass classifiers of ELM.

- (1) *Multiclass classifier with single output*: ELM can approximate any target continuous functions and the output of the ELM classifier $\mathbf{h}(\mathbf{x})\beta$ can be as close to the class labels in the corresponding regions as possible. Thus the classification problem for ELM with a single-output node can be formulated as [21]:

$$\begin{aligned} \text{Minimize : } L_{\text{ELM}} &= \frac{1}{2} \|\beta\|^2 + \frac{C}{2} \sum_{i=1}^n \varepsilon_i^2 \\ \text{Subject to : } &\mathbf{h}(\mathbf{x}_i)\beta = t_i - \varepsilon_i, \quad i = 1, \dots, n. \end{aligned} \quad (4)$$

For multiclass problems, among all the multiclass labels, the predicted class label of a given testing sample is the closest to the output of ELM classifier.

- (2) *Multiclass classifier with multioutputs*: If ELM has multioutput nodes, an m -class classifier is corresponding to m output nodes. If the original class label is l , the expected output vector

of the m output nodes is $\mathbf{t}_l = \overbrace{[0, \dots, 0, 1, 0, \dots, 0]}^l$. That is, the l th element of $\mathbf{t}_l = [t_{l,1}, \dots, t_{l,m}]^T$ is one and the rest of the elements are zero. The classification problem for ELM with multioutput nodes is [21]:

$$\begin{aligned} \text{Minimize : } L_{\text{ELM}} &= \frac{1}{2} \|\beta\|^2 + \frac{C}{2} \sum_{i=1}^n \|\varepsilon_i\|^2 \\ \text{Subject to : } &\mathbf{h}(\mathbf{x}_i)\beta = \mathbf{t}_i^T - \varepsilon_i^T, \quad i = 1, \dots, n. \end{aligned} \quad (5)$$

where $\varepsilon_i = [\varepsilon_{i1}, \dots, \varepsilon_{im}]^T$ is the training error vector of the m output nodes with respect to the training sample \mathbf{x}_i .

If a feature mapping $\mathbf{h}(\mathbf{x})$ is unknown to users, the output function of ELM classifier is

$$f(\mathbf{x}) = \mathbf{h}(\mathbf{x})\mathbf{H}^T \left(\frac{\mathbf{I}}{C} + \mathbf{H}\mathbf{H}^T \right)^{-1} \mathbf{T} = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n)] \left(\frac{\mathbf{I}}{C} + \mathbf{M} \right)^{-1} \mathbf{T}, \quad (6)$$

where $\mathbf{M} = \mathbf{H}\mathbf{H}^T$, $m_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ and $k(\mathbf{x}, \mathbf{y})$ is a positive semi-definite kernel function. If a feature mapping $\mathbf{h}(\mathbf{x})$ is known, we have $\mathbf{h}(\mathbf{x}) = [G(a_1, b_1, \mathbf{x}), \dots, G(a_L, b_L, \mathbf{x})]$, where $G(a, b, \mathbf{x})$ is a non-linear piecewise continuous function satisfying ELM universal approximation capability theorems [22–26] and $\{(a_i, b_i)\}_{i=1}^L$ are randomly generated according to any continuous probability distribution. The output function of ELM classifier is

$$f(\mathbf{x}) = \mathbf{h}(\mathbf{x})\mathbf{H}^T \left(\frac{\mathbf{I}}{C} + \mathbf{H}\mathbf{H}^T \right)^{-1} \mathbf{T}, \quad (7)$$

or

$$f(\mathbf{x}) = \mathbf{h}(\mathbf{x}) \left(\frac{\mathbf{I}}{C} + \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{T}, \quad (8)$$

Download English Version:

<https://daneshyari.com/en/article/407715>

Download Persian Version:

<https://daneshyari.com/article/407715>

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