

# PILAE: A Non-gradient Descent Learning Scheme for Deep Feedforward Neural Networks

P. Guo, *Senior Member, IEEE*, X. L. Zhou, and K. Wang

**Abstract**—In this work, a non-gradient descent learning scheme is proposed for deep feedforward neural networks (DNN). As we known, autoencoder can be used as the building blocks of the multi-layer perceptron (MLP) deep neural network. So, the MLP will be taken as an example to illustrate the proposed scheme of pseudoinverse learning algorithm for autoencoder (PILAE) training. The PILAE with low rank approximation is a non-gradient based learning algorithm, and the encoder weight matrix is set to be the low rank approximation of the pseudoinverse of the input matrix, while the decoder weight matrix is calculated by the pseudoinverse learning algorithm. It is worth to note that only few network structure hyperparameters need to be tuned. Hence, the proposed algorithm can be regarded as a quasi-automated training algorithm which can be utilized in autonomous machine learning research field. The experimental results show that the proposed learning scheme for DNN can achieve better performance on considering the tradeoff between training efficiency and classification accuracy.

**Index Terms**—Autoencoder, Learning scheme, Pseudoinverse learning algorithm, Low rank approximation, Feedforward deep neural network.

## I. INTRODUCTION

Representation learning is to make the classifiers easier perform the given task by learning useful information and extracting essential features from the data [1]. As a way of hierarchical representation learning, deep learning has achieved state-of-the-art performance in a broad range of domains, such as computer vision [2],[3],[4] and speech recognition [5],[6],[7]. There are two paradigms developed parallel in deep learning: one roots in probabilistic graphical models and the other roots

in neural networks [1],[8]. In literature, both of them focus on single layer greedy learning modules and their variants. Among the building block modules, two main modules are explored intensively. One is the restricted Boltzmann Machine (RBM) that is the representation of the probabilistic model, while the network architecture is the same with multilayer neural network in practical implementation [35]. The other is the autoencoder that is the representation of the traditional multi-layer perceptron (MLP).

Autoencoder [9],[10] is a simple artificial neural network that utilizes the input as the output to learn the representation of the raw data in an unsupervised fashion. Thus, autoencoder can be viewed as an unsupervised feature learning model. Several variants of autoencoder are developed, such as sparse autoencoders for overcomplete representation [11], denoising autoencoders and contractive autoencoders for learning robust representation [12],[13]. Autoencoder has widely applications, for example, constrained autoencoders can be used for enhanced understanding of data [40]. Autoencoders can be used as the building blocks for constructing deep networks. During the construction process, deep neural network is trained in a greedy layer-wise scheme, and then the trained autoencoders (without the decoder) are stacked into a multi-layer neural model. In the stacked autoencoders, the output of the previous autoencoder is used as the input of the subsequent one. The output of different hidden layer of the network can be views as the different level features of the original data.

Most of the learning algorithms for training autoencoders are based on the gradient descent algorithm, such as the back propagation (BP) algorithm. The gradient descent-based algorithm is an iterative optimization algorithm which is time

This work is fully supported by the grants from the National Natural Science Foundation of China under Grant No. 61375045, the Joint Research Fund in Astronomy (Grant No.U1531242) under cooperative agreement between the National Natural Science Foundation of China (NSFC) and Chinese Academy of Sciences (CAS), and Beijing Natural Science Foundation under Grant No. 4162027.

P. Guo is with the School of Systems Science, Beijing Normal University, Beijing, China, 100875 (e-mail: pguo@bnu.edu.cn). He is the author to whom all the correspondence should be addressed.

X. L. Zhou is with the Department of Technology and Industry Development, Beijing City University, Beijing, China, 100083 (e-mail: zxlmouse@bcu.edu.cn).

K. Wang is with the Department of Electrical and Computer Engineering, University of Michigan-Dearborn, Dearborn, MI, USA, 48128 (e-mail: wangke@bit.edu.cn).

consuming. Hence, the main problem associated with stacked autoencoders based deep learning is the high computation cost, especially when the data set is large and the dimension of input data is high. Besides, most gradient descent-based algorithms require carefully tuning of several control parameters. These hyperparameters, such as maximum epoch, step length and momentum, are crucial to the success of the algorithm. However, the hyperparameters are usually related to specific application. Hence it is dependent on empirical tricks to select the appropriate hyperparameters. Furthermore, the design of the network structure, i.e. the number of the hidden layers and the number of neurons in each hidden layer, is also a very difficult task, and it is made by trial and error in general.

The connection weights in deep neural network are significantly redundant [14]. Most of the weights can be predicted by using a small fraction of them. In order to improve the training efficiency, it is proposed to reduce the weights size during training by pruning any weight whose value approaches zero [15]. Re-parameterization is an alternative way to reduce the weights size, which is the transformation of parameterization space of a model [16]. Constraining the singular values of weight matrix [17], weight normalization [18] and factorizing the weight matrix by periodic functions [19] can accelerate the training speed to some extent. However, the iterative gradient decent based learning algorithm adopted in those methods dominates the time consuming.

From literature review, we know that the gradient descent-based algorithms either fail or suffer from significant difficulties for some types of simple problems in deep learning [20]. It is necessary to develop non-gradient based learning algorithms for deep learning. In this work, a fast and quasi-automated learning algorithm, that is, pseudoinverse learning algorithm with low rank approximation is proposed for training the autoencoder which will be used as the building blocks of the deep neural network. The proposed algorithm is based on the pseudoinverse learning algorithm which can train the feedforward network in an analytical way [21], [22], [24]. During the training process, the rank and the low rank approximation [26] of the input matrix are obtained by singular value decomposition (SVD) technique. The rank of the input matrix can be used to guide the setting of the number of hidden neurons. The encoder weight matrix is set to be the low rank approximation of the pseudoinverse of the input matrix. The decoder weight matrix is calculated by the pseudoinverse learning algorithm. The depth of the DNN is determined by the criterion that the output matrix of the hidden layer is close to a full rank matrix. It is worth to note that there is no hyperparameters in training algorithm need to be tuned, and suggestion is given for the network structure related hyperparameter selection. Hence, the proposed algorithm can be regarded as a quasi-automated training algorithm for deep feedforward neural networks.

## II. RELATED WORKS

### A. Multi-layer Perceptron

Multi-layer perceptron is a kind of feedforward neural networks, which was most studied in mid-eighties of last century. In our previous work [23], when error distribution is assumed as Gaussian, under the framework of the Kullback–Leibler (KL) divergence, the error (loss) function has following form:

$$J \approx \frac{1}{2N\sigma^2} \sum_{i=1}^N (\|z_i - g(x_i, W)\|^2 + h_x \|g'(x_i, W)\|^2), \quad (1)$$

where  $g(x_i, W)$  is a neural network mapping function,  $x_i$  is the input signal vector,  $z_i$  is the corresponding target output vector, and  $h_x$  is the smooth regularization parameter. When loss function with the form of Eq. (1) is applied to regularized autoencoder, it is the same with empirical loss in Ref. [41].

If we consider linear function mapping approximation in regularization term only, i.e.  $g(x, W) = Wx$ , the regularization term becomes a weight decay form as follows.

$$J \approx \frac{1}{2N\sigma^2} \sum_{i=1}^N \|z_i - g(x_i, W)\|^2 + \frac{1}{2\sigma^2} h_x \sum_{j=1}^M w_j^2, \quad (2)$$

where  $M$  is the number of network weight parameters and  $w_j$  is an element of the matrix  $W$  in a vector expression.

The regularization parameter  $h_x$  can be estimated based on training data as [23]:

$$h_x \approx d_x^2 [1 + (d_x - 1)^2] \frac{\sum_{i=1}^N \|z_i - g(x_i, W)\|^2}{N \sum_{j=1}^M w_j^2}, \quad (3)$$

where  $d_x$  is the dimension of input vector  $x$ .

### B. Autoencoder

An autoencoder can be viewed as a simple feedforward neural network with one input layer, one hidden layer, and one output layer [9],[10]. The goal of training an autoencoder is to learn data representation with minimal reconstruction error.

From the input layer to the hidden layer, the encoder represents the input vector  $x$  as a feature vector  $y$ , which is realized by linear mapping and nonlinear activation function:

$$y = f(W_e x). \quad (4)$$

From the hidden layer to the output layer, the decoder reconstructs the input vector  $x$  as vector  $z$ :

$$z = g_1(W_d y). \quad (5)$$

In Eq. (4) and Eq. (5),  $W_e$  and  $W_d$  are the connection weight matrices,  $f$  and  $g_1$  are activation functions. For concision, the bias parameters are omitted and  $g_1$  is set to be linear function in this paper. If the bias neuron is considered, it is simple to expand weight matrix as an augmented matrix [30]. Weight parameters in an autoencoder are estimated by minimizing the average reconstruction error function:

$$J = \frac{1}{N} \sum_{i=1}^N \|x_i - W_d(f(W_e x_i))\|_2^2, \quad (6)$$

where  $N$  is the number of input data. Error function in Eq. (6) is the sum-of-square error function. If we study regularized autoencoder, empirical loss function in Eq. (1) should be considered.

The weights  $W_d$  can be set as  $W_d = W_e^T$ , which is termed as tied weights so as to decrease the number of estimated free parameters in an autoencoder. Autoencoder can be seen as a representation transformation. When the number of hidden neurons is constrained to be less than the number of input neurons, a compressed data representation of the input data is obtained, which realizes the dimension reduction (feature extraction) task. The constraints of tied weights and hidden neuron number (information bottleneck) can prevent the autoencoder from learning identity mapping.

### C. Basic Pseudoinverse Learning Algorithm

PseudoInverse Learning (PIL) Algorithm is a fast learning algorithm for feedforward neural networks proposed by Guo *et al* [21],[22],[24]. The basic idea of the PIL algorithm is to find a set of orthogonal vector bases and use the nonlinear activation function to make the output vector of the hidden layer neurons orthogonal, and then obtain the output weights of the network by calculating the pseudoinverse solution.

Suppose we have a given data set denoted as  $D = \{x_i, o_i\}_{i=1}^N$ , where  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})^T \in R^d$  is the input vector and  $o_i = (o_{i1}, o_{i2}, \dots, o_{im})^T \in R^m$  is the desired output vector. Denote  $X \in R^{d \times N}$  as the input matrix and  $O \in R^{m \times N}$  as the desired output matrix. (If input bias neuron is considered, input matrix  $X$  will be  $X \in R^{(d+1) \times N}$ ). Consider a feedforward neural network with  $L$  layer, where the linear activation function is used in last layer. Let  $W^l$  be the weight matrix which connects the layer  $l$  and the layer  $l+1$ . Denote  $H^l = f(W^{l-1}H^{l-1})$  as the output matrix of the layer  $l$  where  $f(\mathbf{x})$  is an activation function.

The task of training the network is trying to find the optimal weight matrices which minimize the error function, for example, most used sum-of-square error function:

$$\min \|W^L H^L - O\|^2. \quad (7)$$

Based on the generalized linear algebra theory [27], for the

equation  $W^L H^L = O$ , the optimal approximation solution is  $W^L = O(H^L)^+$ , where  $(H^L)^+$  represents the pseudoinverse of the matrix  $H^L$ . Let  $W^L = O(H^L)^+$ , the objective function becomes:

$$\min \|W^L H^L - O\|^2 = \|O(H^L)^+ H^L - O\|^2. \quad (8)$$

For solving the above optimization problem, the connection weight  $W^L$  is set to be the pseudoinverse matrix  $(H^L)^+$ , the output matrix  $H^L$  is propagated feedforward and the rank of  $H^L$  is increased gradually by the nonlinear activation function transformation. Once the rank of  $H^L$  is close to full, that is,  $(H^L)^+ H^L$  is close to an identity matrix, namely, if  $\|(H^L)^+ H^L - I\|^2 < \varepsilon$  the learning procedure can be terminated. This can be considered as a set of orthogonal vector bases is reached with nonlinear activation of hidden neurons,  $(H^L)^+ H^L$  is the orthogonal projector which project output matrix to approximate its self space spanned by the column vectors.

The PIL algorithm is a feedforward algorithm without any iterative optimization process. Meanwhile, it is also an automated algorithm without critical user-specified parameters such as learning rate or momentum constant. Thus, PIL algorithm is more efficient than the back propagation and other gradient descent-based algorithm. Furthermore, the structure of the network generated by PIL algorithm is the one with deep attribution, and the dynamical growth in depth is dependent on the given training data set.

### D. The Low Rank Approximation

Low rank matrix approximation is well studied in the numerical linear algebra community [26]. Many classical matrix decomposition techniques are adopted for low rank approximation, such as singular value decomposition (SVD) and QR decomposition [28]. In our proposed method, truncated SVD is used to calculate the rank of matrix and the low rank approximation of the pseudoinverse matrix.

For any matrix  $A \in R^{m \times n}$ , it can be factorized by SVD as:

$$A = U \Sigma V^T, \quad (9)$$

where  $U \in R^{m \times m}$  is an orthogonal matrix whose columns are the eigenvectors of  $AA^T$ ,  $V \in R^{n \times n}$  is an orthogonal matrix whose columns are the eigenvectors of  $A^T A$ , and  $\Sigma \in R^{m \times n}$  is an diagonal matrix whose entries are in descending order,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$ , along the main diagonal:

$$\Sigma = \text{diag}\{\sigma_1, \dots, \sigma_r, 0, \dots, 0\}, \quad (10)$$

with  $r = \text{rank}(\mathbf{A})$ . In the above equation,  $\sigma_1, \dots, \sigma_r$  are the square roots of the eigenvalues of  $\mathbf{A}^T \mathbf{A}$ , which are called the singular values of  $\mathbf{A}$ .

Mathematically, truncated SVD applied to a matrix will produce a low rank approximation to that matrix.

### III. PROPOSED FAST LEARNING SCHEME FOR MULTILAYER DEEP NEURAL NETWORK

For a given task with data set,  $D = \{\mathbf{x}_i, \mathbf{z}_i\}_{i=1}^N$ , we assume that the network mapping function  $g(\mathbf{x}, \mathbf{W})$  in Eq. (1) is a deep feedforward neural network (DNN), which has one input layer, one output layer and several hidden layers.

$$g(\mathbf{x}, \mathbf{W}) = \mathbf{W}^l f(\mathbf{W}^{l-1} \dots f(\mathbf{W}^1 f(\mathbf{W}^0 \mathbf{x})) \dots). \quad (11)$$

In our proposed learning scheme, “divide and conquer” strategy is utilized to train the DNN. With weight decay regularization loss function in Eq. (2), the formal solution of  $\mathbf{W}^l$ , which connects last hidden layer to output layer, is the pseudoinverse solution:

$$\mathbf{W}^l = \mathbf{Z}\mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T + \lambda \mathbf{I})^{-1} = \mathbf{Z}\mathbf{Y}^+. \quad (12)$$

Where  $\mathbf{Z}$  is target output matrix,  $\mathbf{Y}$  is the last hidden layer output matrix,

$$\mathbf{Y} = f(\mathbf{W}^{l-1} \dots f(\mathbf{W}^1 f(\mathbf{W}^0 \mathbf{X})) \dots), \quad (13)$$

and  $\mathbf{Y}^+$  has the form  $\mathbf{Y}^+ = \mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T + \lambda \mathbf{I})^{-1}$ . This form is similar with limit form of Moore-Penrose inverse, but here  $\lambda$  is the regularization parameter and it should not approach to zero. The weight matrices in Eq. (13) will be trained in an unsupervised layered-wise manner, that is, the network in Eq. (13) is taken as a stacked autoencoders as shown in Fig.1. Following we will describe pseudoinverse learning algorithm for autoencoders.

#### A. Proposed Pseudoinverse Learning Algorithm with Low Rank Approximation for Autoencoders

For a data matrix  $\mathbf{X} \in R^{d \times N}$ ,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ , where  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T \in R^d$  is the  $i$ -th sample with dimension  $d$ . For training autoencoder, the objective (loss) function (reconstruction error) can be written as:

$$\min \|\mathbf{W}_d \mathbf{H} - \mathbf{X}\|^2. \quad (14)$$

The pseudoinverse approximate solution is used to solve this optimization problem:

$$\mathbf{W}_d = \mathbf{X}\mathbf{H}^+. \quad (15)$$

According to the basic PIL algorithm, the connection weight  $\mathbf{W}_e$  should be set as  $\mathbf{X}^+$ , which means implicitly a constraint that the number of hidden neurons is equal to the number of input samples. In general, the number of input samples is far more than the dimension of the input data.

Here, autoencoder is used to learn data representation. In order to learn good representation of the data, constraints such as sparse and/or information bottleneck (IB) are most used. The constraint of IB is that the number of hidden neurons should be less than the dimension of the input data. With either sparse or IB constraint, the identity mapping can be avoided, while data representation is learnt [1]. According to the manifold hypothesis that real world data presented in high dimensional spaces is likely to concentrate in the vicinity of non-linear sub-manifolds of much lower dimensionality [29], the number of hidden neurons should be equal or more than the rank of the input data matrix.

In order to learn the data representation, a pseudoinverse learning algorithm with low rank approximation for autoencoders is proposed in this paper. In the proposed method, the rank of the input matrix is used to guide the setting of the hidden neuron number. The encoder weight matrix is set to be the low rank approximation of the pseudoinverse of the input matrix, and the decoder weight matrix is calculated with PIL algorithm [25].

There are three merits for using low rank approximation with SVD. First, the number of hidden neurons can be set automatically with a formula according to the rank of input data matrix. Second, the encoder weight matrix can be calculated with truncated SVD. Third, those samples, which can be expressed as the linear combination of other samples, are merged. By merging the dependent samples by low rank approximation and encoding the input data by linear dimension reduction, the input data are mapped into a feature space, in which its intrinsic characteristics can be better reflected.

The implementation of the pseudoinverse learning algorithm with low rank approximation is described in detail as follows.

1) Determine the number of hidden neurons by a formula with the rank of input data matrix.

The rank of input matrix can be calculated by the SVD method. At first, the input data matrix  $\mathbf{X}$  can be decomposed by the SVD as follows:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (16)$$

where  $\mathbf{U} \in R^{d \times d}$  and  $\mathbf{V} \in R^{N \times N}$  are orthogonal matrices,  $\mathbf{\Sigma} \in R^{d \times N}$  is a diagonal matrix where the diagonal elements are the singular values of  $\mathbf{X}$ . Then the rank of the input matrix  $\mathbf{X}$  is obtained by counting the number of nonzero singular values in matrix  $\mathbf{\Sigma}$ :

$$r = \text{rank}(\mathbf{X}) = \text{the number of nonzero}(\Sigma). \quad (17)$$

The number of hidden neurons is set according to the rank of input matrix  $\mathbf{X}$ . In order to avoid obtaining an identity mapping in the autoencoder, the number of hidden neurons should generally be less than the dimension of the input data. However, too less number of hidden neurons will lead to more loss of information during the process of mapping the input data to the hidden feature space, which will inevitably result in higher reconstruction error. Hence, the number of hidden neurons should not be too small. In our previous work [25], on considering the tradeoff between feature extraction and reconstruction error, the number of hidden neurons is assigned as

$$p = r + \alpha(d - r), \alpha \in [0, 1], \quad (18)$$

where  $d$  is the dimension of the input data  $\mathbf{x}$ . In the case of the data matrix with *i.i.d.* noise, whose rank is almost full, we can apply data preprocessing techniques to reduce the noise, or we can force the dimensionality of input data to reduce. Therefore, the number of hidden neurons can be set with a decay factor shown as follows:

$$p = \beta d, \beta \in (0, 1]. \quad (19)$$

2) Calculate the encoder weight matrix as the low rank approximation of the pseudoinverse of the input matrix. The parameter  $\alpha$  or  $\beta$  is estimated based on reconstruction error.

From the result of the SVD of  $\mathbf{X}$ , the pseudoinverse of matrix  $\mathbf{X}$  can be calculated as follows:

$$\mathbf{X}^+ = \mathbf{V}\Sigma^+\mathbf{U}^T, \quad (20)$$

where  $\Sigma^+ \in R^{N \times d}$  is the diagonal matrix composed of the reciprocal of nonzero elements in matrix  $\Sigma$ .

Let

$$\hat{\mathbf{X}}^+ = \hat{\mathbf{V}}\Sigma^+\mathbf{U}^T \in R^{p \times d}, \quad (21)$$

be the  $p$  low rank approximation of  $\mathbf{X}^+$ , where  $\hat{\mathbf{V}} \in R^{p \times N}$  is the matrix composed of the first  $p$  rows of  $\mathbf{V}$ . This approximation is similar with truncated SVD method, because the truncated SVD can be viewed as a kind of regularization [36], hence it can serve as a dimension reduction technique. In mathematics, this also can be realized by defining a  $p \times N$  diagonal matrix  $\Omega$ , in which those diagonal elements are 1 and non-diagonal elements are 0, and  $\hat{\mathbf{V}} = \Omega\mathbf{V}$  [42]. The encoder weight matrix is set to be the low rank approximation of the pseudoinverse matrix:

$$\mathbf{W}_e = \hat{\mathbf{X}}^+. \quad (22)$$

Here, the encoder matrix is calculated with the truncated SVD method. Then the input matrix  $\mathbf{X}$  is mapped into the  $p$ -

dimensional hidden feature space as

$$\mathbf{H} = f(\mathbf{W}_e\mathbf{X}) \in R^{p \times N}. \quad (23)$$

Here,  $f(\cdot)$  is the activation function which is usually a nonlinear function, for example, the sigmoid function:

$$f(x) = \frac{1}{1 + \exp(-x)}, \quad (24)$$

or the hyperbolic function

$$f(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}, \quad (25)$$

or step function

$$f(x) = \begin{cases} 1 & \text{if } x \geq 0.5 \\ 0 & \text{if } x < 0.5 \end{cases}. \quad (26)$$

and so on [30].

3) Calculate the decoder weight matrix with the PIL algorithm.

According to the basic idea of PIL algorithm, the optimal solution of equation of  $\mathbf{W}_d\mathbf{H} = \mathbf{X}$  is the pseudoinverse solution

$\mathbf{W}_d = \mathbf{X}\mathbf{H}^+$ . Hence, the pseudoinverse of  $\mathbf{H}$  needed to be calculated. It can also be calculated by using SVD or QR method. Alternatively, if  $\mathbf{H}$  is a row full rank matrix, it can be calculated as

$$\mathbf{H}^+ = \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}. \quad (27)$$

When  $\mathbf{H}$  is a rank defective matrix. To avoid ill-posed problem, a weight decay regularization term is added to the objective function as in Eq. (2). This form is also called ridge regression in statistics, and when omitting some optimization irrelative constants, the weight decay regularization in the matrix form is shown as follows:

$$J = \frac{1}{2} \|\mathbf{W}_d\mathbf{H} - \mathbf{X}\|^2 + \frac{\lambda_1}{2} \|\mathbf{W}_d\|^2. \quad (28)$$

The regularization solution to the optimization problem in above equation has analytical solution. Taking derivative of  $J$  to  $\mathbf{W}_d$  and let it be zero, it can be obtained

$$(\mathbf{W}_d\mathbf{H} - \mathbf{X})\mathbf{H}^T + \lambda_1\mathbf{W}_d = 0. \quad (29)$$

Then the pseudoinverse of  $\mathbf{H}$  is calculated as follows:

$$\mathbf{H}^+ = \mathbf{H}^T(\mathbf{H}\mathbf{H}^T + \lambda_1\mathbf{I})^{-1}, \quad (30)$$

where  $\lambda_1$  is a regularization parameter.

The decoder weights can be calculated as:

$$W_d = XH^T(HH^T + \lambda_1 I)^{-1} = XH^+ \quad (31)$$

Here, the regularization parameter  $\lambda_1$  can be regarded as a smooth parameter as in Eq. (3), so it can be estimated with training data only. This regularization parameter controls the reconstruction error, here we set it to be a very small value to guarantee that matrix  $(HH^T + \lambda_1 I)$  is not singular.

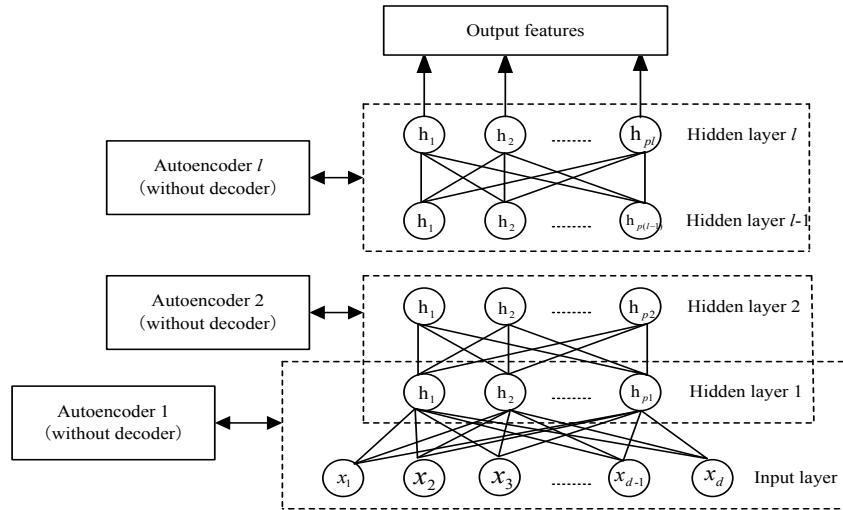
After the decoder weight is obtained, the decoder weights and the encoder weights will be tied, that is,  $W_e = W_d^T$ . The tied weights can reduce the freedom degree of the model. The hidden feature matrix  $H = f(W_e X)$  can be used as the input data matrix for succeeding autoencoders.

### B. Pseudoinverse Learning for Deep Stacked Autoencoders

Based on the proposed pseudoinverse learning algorithm with low rank approximation for autoencoders, a fast and quasi-automated learning algorithm for stacked autoencoders is developed, which can overcome the shortcomings of the gradient descent-based training algorithms for deep feedforward neural network such as BP algorithm.

In the proposed method, the pseudoinverse learning algorithm with low rank approximation is applied for training the autoencoder which will be used as the building blocks of the stacked autoencoders. After obtaining the encoder connection weight matrix  $W_e$ , the hidden output matrix  $H = f(W_e X)$  is used as the input matrix for succeeding autoencoders. This procedure is continued until a criterion is met. This criterion making the network stop growing may adopt a predesigned maximum number of layers, or utilize the  $\|H^+H - I\|^2 < \varepsilon$  as we did in previous work [24].

In fact, the proposed learning scheme can be explained as the greedy layer-wise learning strategy applied to the DNN, in which all trained autoencoders are then stacked to construct a deep neural network by removing the decoder. For the more complex networks, the number of the weight parameters is reduced by tied weight operation, which means taking the transposed decoder weight matrix as the encoder weight matrix. The structure of the deep neural network based on stacked pseudoinverse autoencoders is shown in Fig. 1, while the mathematical expression is in Eq. (13).



**Fig. 1** Structure of the stacked pseudoinverse learning autoencoders deep network

After training, the output feature matrix of the stacked autoencoders based deep neural network can be calculated as Eq. (13), and the weight matrices are

$$W^i = W_e^i, \quad i = 0, 1, \dots, l-1. \quad (32)$$

With learnt feature matrix, we can use Eq. (12) to calculate the output weight matrix. If users intend to use the features for classification task and do not utilize single hidden layer neural net (SHLN) as we did, an additional classifier should be used, such as softmax, support vector machine or other relative simple classifiers, to obtain final results.

In order to realize global optimization for the DNN structure, in practice, we need to slightly adjust the neuron number  $p_L$  and output layer regularization parameter  $\lambda$  of the last hidden layer to minimize generalization error. Here we propose an empirical formula to estimate  $p_L$  based on the rank of input data and the number of training samples:

$$p_L = \begin{cases} P(r, N) & \text{if } P(r, N) > 0 \\ r + \alpha(d - r) & \text{otherwise} \end{cases} \quad (33)$$

where

$$P(r, N) = \lceil \theta_0 + \theta_1 r + \theta_2 N + \theta_3 r^2 + \theta_4 N^2 \rceil, \quad (34)$$

and  $r$  is the rank of the input data,  $N$  represents the number of training samples,  $\theta_i$  ( $i = 1, 2, 3, 4$ ) are constant coefficients.

#### IV. EXPERIMENT RESULTS

In this section, we first conduct experiment 1 to demonstrate how to determine the coefficients in Eq. (34). Then experiment 2 is conducted to compare our proposed method with the baseline. Since the BP algorithm is the representative gradient descent-based algorithm for training the autoencoders, the proposed algorithm is compared with the BP algorithm on different data sets to evaluate its performance. The deep neural network is trained in the classical greedy layer-wise scheme by the two algorithms respectively. Each hidden layer is associated with a trained autoencoder (without the decoder). The output of the last hidden layer can be used as the input features for specific learning tasks, e.g. classification.

##### A. Experiments Setup

In the first experiment, we train SHLNs with several publicly available data sets. These data sets are different in terms of feature number, rank, and sample number. For each data set, the relationship between the hidden neuron number and the classification performance is analyzed. Specifically, a series of SHLNs with different hidden neuron number are trained on each data set so as to find the optimal network structure which can achieve the highest test accuracy for each data set. 5-fold cross validation is employed for evaluation of all networks. With the obtained optimal network structure, regression analysis is then used to find the relationship between the rank of data along with the sample number and the optimal hidden neuron number, i.e. the equation defined in Eq. (33).

In the second experiment, the proposed algorithm and the BP algorithm are used to train networks with the same structure, respectively. After data representation is learnt, classification task is performed. To be specific, the output of the stacked autoencoders is taken as the input features of the classifiers, such as softmax and SHLN. The softmax model is a generalization of logistic regression for multi-class problems.

In the softmax model, the probability of the given sample belonging to each class is estimated. And then, the class label of the given sample is determined by the maximum posterior probability. For a given sample  $\mathbf{x}_i$ , the output of the softmax is the estimated posterior probability calculated as:

$$p(y_i = k | \mathbf{x}_i; \Theta) = \frac{\exp(\theta_k^T \mathbf{x}_i)}{\sum_{j=1}^K \exp(\theta_j^T \mathbf{x}_i)}, \quad (35)$$

where  $y$  is the class label and  $k \in \{1, 2, \dots, K\}$ ;  $\Theta$  is the parameter set of the classifier.

Given a training set, the parameter set  $\Theta$  is estimated by minimizing the following loss function:

$$J(\Theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K 1_{\{y_i=j\}} \log p(y_i = j | \mathbf{x}_i, \Theta), \quad (36)$$

where  $1_{\{y_i=j\}}$  is an indicative function:

$$1_{\{y_i=j\}} = \begin{cases} 1 & y_i = j \\ 0 & \text{otherwise} \end{cases}. \quad (37)$$

This optimization problem can be solved with a gradient descent-based algorithm. For fairness, the structure of the neural network for BP algorithm is specified as same as that of the proposed algorithm.

When SHLN classifier is adopted as the last layer of neural network, the weights of the last layer can be calculated with Eq. (12).  $g(\mathbf{x}, \mathbf{W})$  expressed in Eq. (11) is used for classification task. If the output of  $g(\mathbf{x}, \mathbf{W})$  is feed into a softmax classifier, the classification results will be converted to probability expressions, and this system is called cascade neural network system.

In the comparison experiment, the network structure of different training algorithms can both be quasi-automatically generated by the proposed learning scheme. Besides, the softmax models share the same parameters in both algorithms. Note that the main purpose of the experiments is not to pay attention to the classification accuracy, but to compare the performance of the feature learning capacity of different algorithms under the same conditions.

TABLE I  
THE SUMMARY OF THE DATA SETS

Data set	#Instances	#Features	Rank
Bupa	250	6	6
Segment	1700	19	19
Mfeat	1400	649	646
Prior	2500	785	621
Har 1	3000	562	540
Advertisement	3279	1558	727
Gina_agnostic	3468	970	970
Spambase	4601	57	57
Isolet	7797	617	617
Har 2	7352	562	540
Sylva	10000	216	204

### B. Data Sets

The benchmark data sets used to analyze the relationship between the rank, sample number and the optimal hidden neuron number are from the OpenML community (<http://openml.org>) [43] and UCI repository [44]. The data sets are described in Table I. Besides, MNIST and Fashion-MNIST are used to compare the performance of the baseline method and the proposed one. MNIST is a common data set which is usually used to evaluate the performance of the deep learning model. MNIST consists of 70,000 handwritten digit images with  $28 \times 28$  pixels in gray scale. The digit images belong to 10 categories: 0-9. Compared with MNIST, Fashion-MNIST [46] poses a more challenging classification task. Being the same with MNIST, it consists of a training set of 60,000 samples and a test set of 10,000 samples. Each sample is a grayscale image of the same size with MNIST. The images belong to 10 categories.

### C. Results and Analysis

#### C.1 Results of experiment 1

To verify the obtained relationship between the rank, the sample number and the optimal hidden neuron number, a leave-one-out cross validation is employed. To be specific, we first leave aside a single data set as a validation set, and use other data sets to generate the relationship specified in Eq. (34) by bivariate regression analysis. Then we use the relationship to predict the hidden neuron number for the validation set, and train a SHLN with the predicted hidden neuron number on the validation set.

The classification accuracy of this SHLN is then compared with the optimal one which is obtained according to the relationship between the hidden neuron number and the classification

performance, which is illustrated in Fig.2. The results of bivariate regression with the rank, the number of samples and the optimal hidden neuron number are shown in Fig.3. The comparison between the classification performance of the predicted network structure and the optimal one is summarized in Table II, in which the columns of #Hidden neurons\* and Test accuracy\* denote the optimal results obtained by using brute force search. It can be observed that the predicted network structure can achieve near-optimal accuracy except the first two data sets whose #features are extremely low. However, from the results of bivariate regression results shown in Fig. 3, it can be observed that using different data sets derives different regression results. Therefore, the data sets could be categorized according to their rank (or #features) and #instances, and then conduct the regression analysis in each category respectively. For example, all data sets with low #features, but large #instances should be analyzed together while the ones with high #features and medium #instances are analyzed together.

#### C.2 Results of experiment 2

The comparison of elapsed training time on MNIST and Fashion-MNIST are shown in Fig. 4 and Fig. 5 respectively. It can be observed that the proposed algorithm is much efficient than the BP algorithm on both datasets. This is due to that the proposed algorithm is a non-gradient method without iterative optimization, whereas the BP algorithm is a gradient based method which requires iterative optimization. Especially when the data size is large and/or the network is deep, the training with BP algorithm is extremely time-consuming.

Since the network structure can be automatically generated in the proposed learning scheme, we do not need to specify the network structure manually. To conduct the comparison, the same network structure is employed by BP algorithm. Besides, the network trained with BP algorithm and the one trained with our proposed scheme use the same classifier i.e. the Softmax regression model or SHLN, sharing same parameters. The classifiers take the learnt features as the input. The classification accuracy is demonstrated in Table III.

Considering that exactly same classifiers are employed, hence better classification results obtained by using our proposed method means that the neural network trained by the proposed algorithm has stronger representation learning capacity.



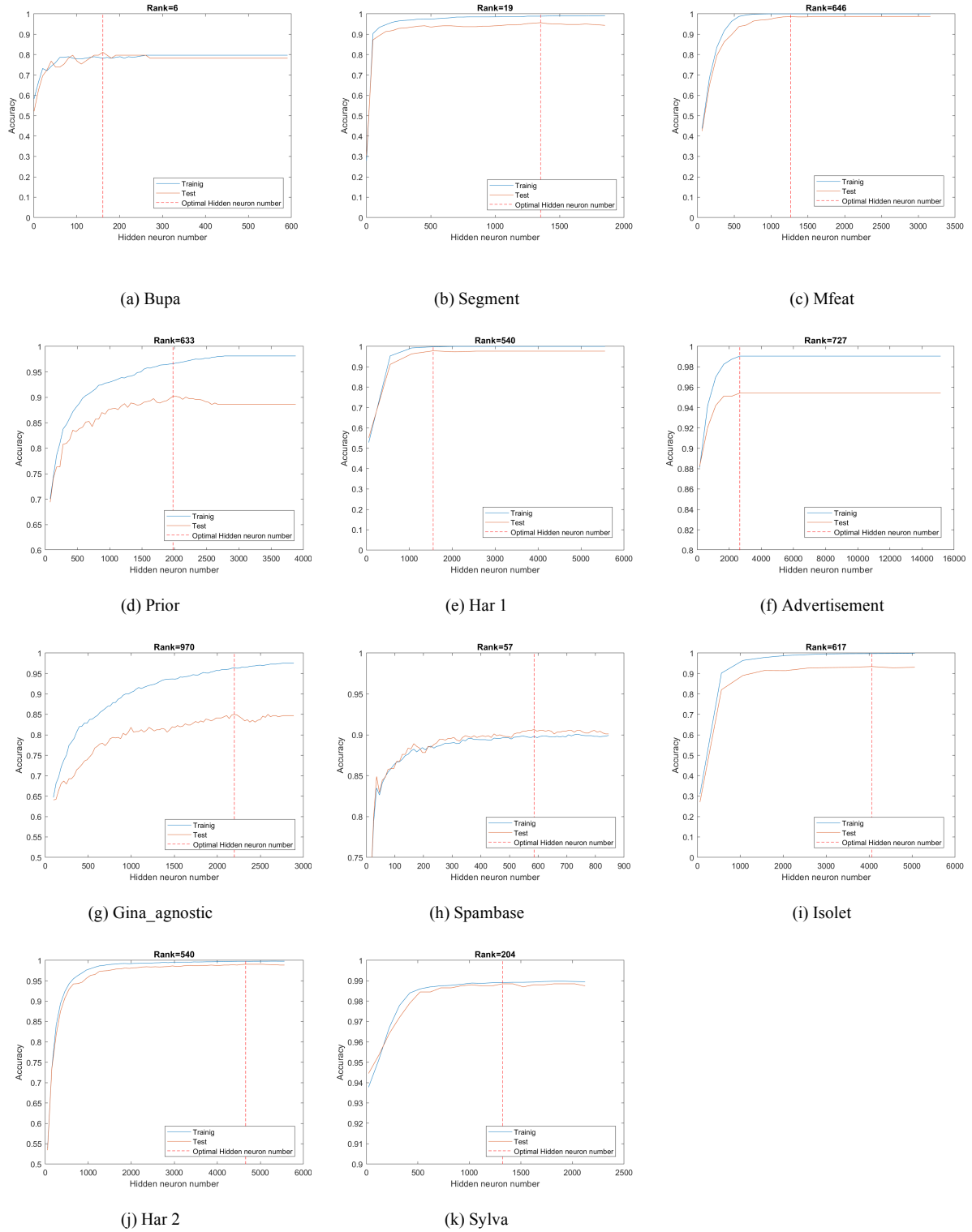


Fig.2: The relationship between the hidden neuron number and the classification performance on different data sets.

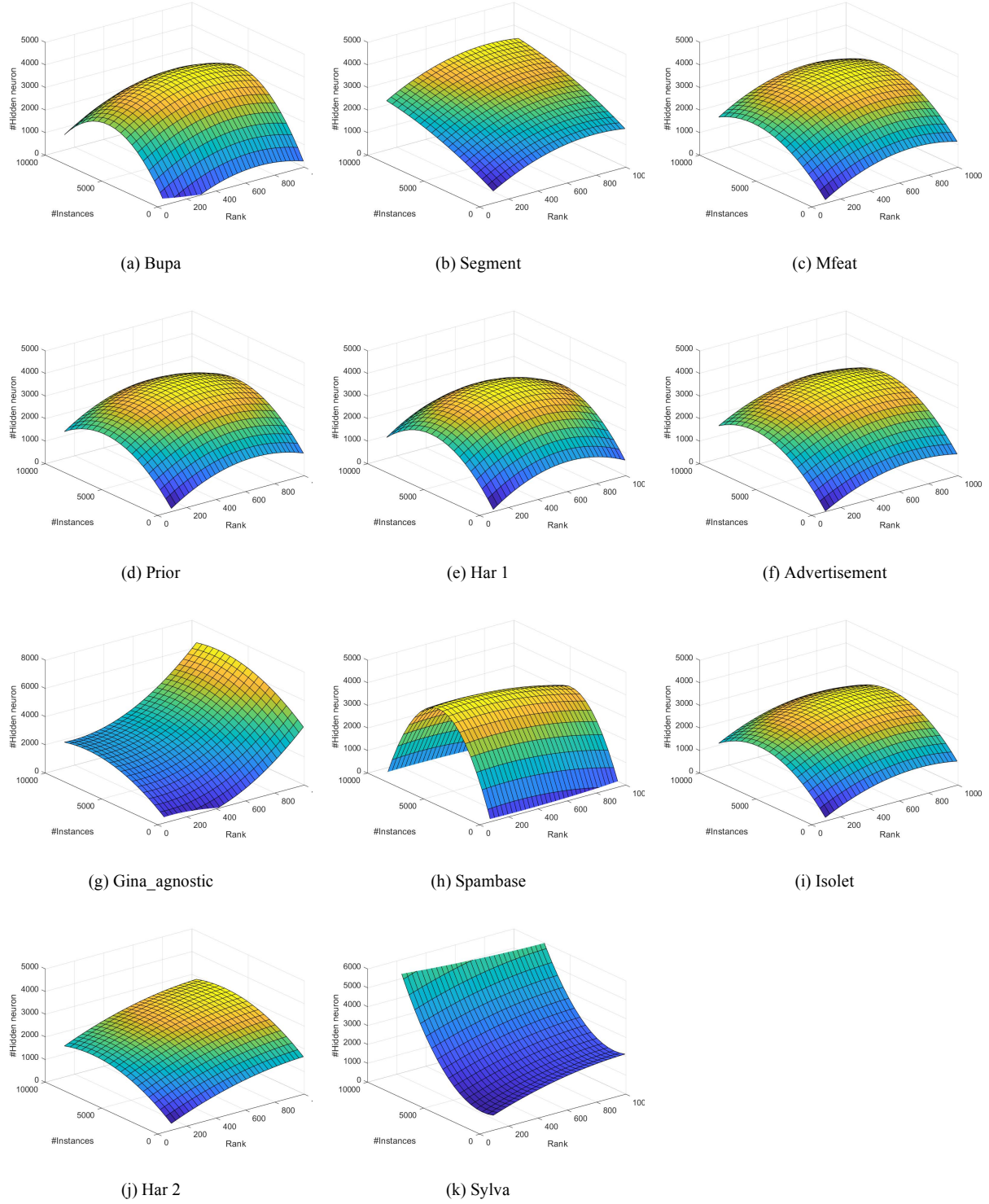


Fig.3: The results of bivariate regression with the rank, the sample number and the optimal hidden neuron number on different by leaving out each data set respectively.

TABLE II  
THE COMPARISON BETWEEN THE PREDICTED NETWORK STRUCTURE AND THE OPTIMAL  
STRUCTURE

Data set	$[\theta_0, \theta_1, \theta_2, \theta_3, \theta_4]$	#Hidden neurons	#Hidden neurons*	Test accuracy	Test accuracy*
Bupa	[-1.2982e3,5.0603,1.0390,-0.0036,-9.0308e-5]	-1064	161	-	81.16
Segment	[-926.2121,5.6153,0.6169,-0.0037,-4.1785e-5]	-52	1352	-	95.67
Mfeat	[-377.4237,4.7423,0.6002,-0.0033-4.5829e-05]	1921	1265	98.75	98.75
Prior	[-434.6510,5.2149,0.6614,-0.0039,-5.4856e-5]	2416	1979	89.32	90.19
Har 1	[-515.3591,5.5879,0.7571,-0.0045,-6.714e-5]	2627	1557	97.67	97.83
Advertisement	[-459.4272,4.2053,0.6779,-0.0029,-5.2552e-5]	2497	2556	95.42	95.42
gina_agnostic	[-373.4045,-1.9230,0.7241,0.0063,-4.7683e-5]	5339	2197	84.70	84.99
spambase	[-712.0285,0.7562,1.8134,-0.0013,-0.0002]	3510	586	91.09	90.65
isolet	[-293.6718,3.8698,0.6184,-0.0026,-5.1742e-5]	2949	4062	92.69	93.39
Har 2	[-58.6140,3.0542,0.4126,-0.0014,-2.9070e-5]	2603	4657	98.03	99.05
Sylva	[487.8985,3.04204,-0.2173,-0.0014,0.0001]	5893	1322	98.61	98.85

TABLE III  
THE COMPARISON OF TEST ACCURACY WITH DIFFERENT  
ALGORITHMS

	BP algorithm		Our proposed	
	Softmax	SHLN	Softmax	SHLN
MNIST	0.9062	0.9692	0.9611	0.9751
Fashion-MNIST	0.7614	0.8545	0.8621	0.8819

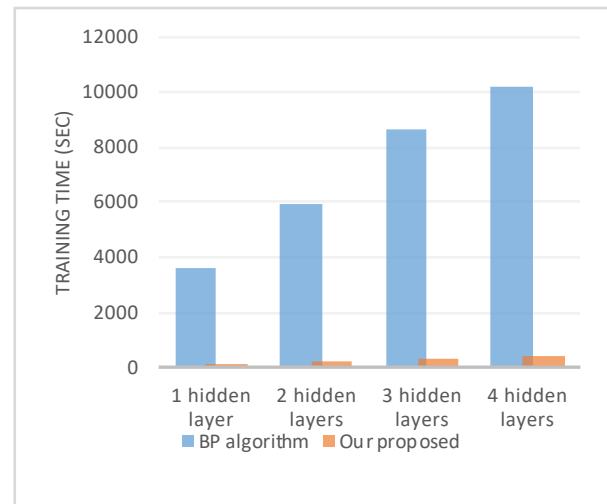


Fig. 4: The comparison of training time of the BP algorithm and our proposed method on MNIST dataset

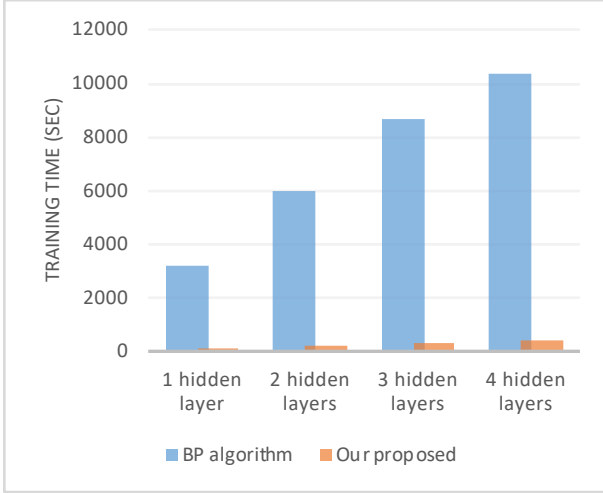


Fig. 5: The comparison of training time of the BP algorithm and our proposed method on Fashion-MNIST dataset

We also analyzed the ratio of rank to dimension in different hidden layers on MNIST and Fashion-MNIST datasets and the results are shown in Fig. 6. For MNIST, after the second hidden layer is added, the rank of the input data equals to the dimension of the input data. For Fashion-MNIST, the input matrix is always row full rank and dimension reduction cannot be achieved by using Eq. (18). Therefore Eq. (19) can be used to reduce the dimension of input matrix. Otherwise, additional denoising methods should be employed to process the input data, such as denoising autoencoder.

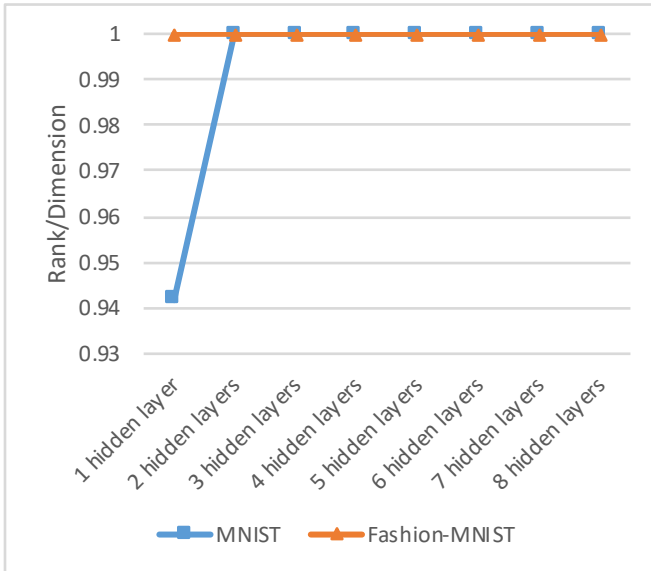


Fig. 6: The ratio of rank to dimension in different hidden layers.

## D. Discussion

### D.1 Advantage

- 1) The proposed learning scheme is a fast learning scheme.

The time complexity analysis of the proposed PIL with low

rank approximation for autoencoders is shown in Table IV. In this table,  $N$  is the number of input data,  $d$  is the dimension of the input data,  $p$  is the number of hidden neurons.

TABLE IV

TIME COMPLEXITY OF THE PIL WITH LOW RANK APPROXIMATION FOR AUTOENCODERS

Operation	Time complexity
SVD(X)	$O(dN^2)$
$\hat{X}^+$	$O(pdN)$
H	$O(pdN)$
$H^+$	$O(p^2N)$
$W_d$	$O(pdN)$
overall	$O(dN^2)$

Generally, there exists  $p < d < N$ , so the overall time complexity of the proposed algorithm for autoencoders is  $O(dN^2)$ . Let  $l$  represents the deep length of the DNN. In practice,  $l$  is far more less than  $p$ . So, the overall time complexity of the proposed algorithm for stacked autoencoders deep neural network is also  $O(dN^2)$ .

Base on above analysis, the proposed algorithm can directly calculate the analytical solution of the optimization object function by using basic matrix operations, such as matrix product and pseudoinverse operation. The connection weight matrices are estimated in a feedforward way without any iterative optimization. Thus, it is a fast learning algorithm. It is more efficiency than the back propagation and other gradient descent-based learning algorithm. This is shown in the above experiments. Also, PILAE learning scheme can work not only in batch learning manner but also in the hybrid of batch and sequence learning manner [45].

- 2) The proposed learning scheme is a quasi-automated learning scheme.

There are two types of hyperparameters involving in the stacked autoencoders deep neural network. One is the model hyperparameters which determine the network's architecture, i.e. the number of the hidden layers (the depth of the network) and the number of hidden neurons in each hidden layer. The other is the control parameter in learning algorithm, such as maximum epoch, step length, weight decay and momentum in BP algorithm. These hyperparameters have an important influence on the performance of the stacked autoencoders deep neural network.

Shown by the procedure of the proposed learning scheme, the weight  $W_e$  connecting the hidden layer and the input layer is estimated by the pseudoinverse of the input matrix with low rank approximation. The weight  $W_d$  connecting the output layer and the hidden layer is the optimal solution minimizing the reconstruction error function. These connection weight matrices are calculated directly in a feedforward way based on generalized linear algebra theory. The regularization parameters  $\lambda_i$   $|_{i=1}^l$  in  $i$ -th hidden layer is fixed as a very small positive value to guarantee that  $H^+$  in Eq. (30) is positive defined. The regularization parameter  $\lambda$  in Eq. (12) can be first estimated with training data as in Eq. (3). And then its optimal estimate is obtained by searching the nearby points around the point estimation. No other control parameters are needed to adjust in the learning process.

Furthermore, the model hyperparameters, that is, the depth of the network and the number of hidden neurons in each hidden layer, are given automatically by the proposed learning scheme. The deep length of the network is determined automatically when the proposed algorithm satisfies the terminal criterion. The hidden layer neuron number in feature learning part is relative to input data rank and data dimension, and determined by an empirical formula, Eq. (18) or Eq. (19). And the parameter  $\alpha$  or  $\beta$  in Eq. (18) or Eq. (19) is estimated based on reconstruction error, that is, with proper reconstruction error setting, we can find a suitable value for  $\alpha$  or  $\beta$ . The hidden layer neuron number in classifier is relative to input data rank and training data number, and estimated by bivariate regression analysis with Eq. (33).

In a word, the proposed learning scheme is not only to improve the learning efficiency comparing the gradient descent-based learning algorithm, but also to provide an empirical designing of the network structure. Thus, the learning speed of the whole neural network system is improved by the proposed learning scheme. Moreover, it is explored preliminarily to realize the autonomous machine learning (AutoML) of the DNN architecture.

3) From the procedure of the proposed algorithm, the encoder weight matrix is tied with decoder weight. As we known,  $W_d = XH^T(HH^T + \lambda I)^{-1}$  which depends on  $W_e$ . If  $W_e$  is changed,  $H$  will be changed also, resulting in  $W_d$  needs to be computed again. In theoretical analysis, it will diverge if weights are iterative computed. This problem might be solved by using batch-normalization techniques.

It is also easy to understand that when encoder weight matrix is calculated as the truncated pseudoinverse of the input data matrix, it will guarantee the values for hidden layer input are almost in the linear region of the sigmoid/Tanh hidden activation function, not in the flat region. This can prevent the saturation of the hidden neuron. It should note that because PILAE is not a gradient descent-based method, it has no

gradient vanish problem.

Our contributions include following points also:

Greedy layer-wise learning strategy is considered as local optima for network's architecture, while in this work global optimization strategy is investigated. For the autoencoder learning, we present an empirical formula to estimate the hidden neuron number. In our learning scheme, searching optimal hyperparameters of DNN structure in multi-dimensional space now is restricted to one dimensional space search problem. In our approaches, only two hyper parameters should be optimized: one is the last hidden layer neuron number, and another is the regularization parameter in system loss function. Also with formula Eq.(3), searching the optimal regularization parameter can be speeded up in a reduced search space. And with Eq. (33), last hidden layer neuron number can be estimated without exhaust search. So our proposed learning scheme will greatly reduce difficulty of the DNN architecture design.

As we known, pre-training deep neural networks with autoencoders is a basic practice of deep learning. However, the learning scheme we presented in this work is not a pre-training algorithm, it is a non-gradient descent learning algorithm for DNN. Layer-wise training the DNN can be regarded as the application of divide and conquer strategy, fine tune is needed for final DNN optimization in traditional pre-training method. Also, as Y. Bengio pointed out in Quora.com, the disadvantage of pre-training deep neural networks with autoencoders is that *"it is greedy, i.e., it does not try to tune the lower layers in a way that will make the work of higher layers easier"*. To avoid too greedy, our strategy is to maintain the reconstruction error in each layer not too small, while search the optimal hidden neuron number in SHLN to reach the better generalization performance of DNN.

In the literature, some researchers estimate latent dimension of the input with SVD method, but relationship between latent dimension and hidden neuron number in autoencoders has not investigated before. In fact, what we presented in Section III.A has illustrated that we do not use SVD to estimate latent dimension of the input, but use truncated SVD to compute the encoder weight. This is corresponding to the low rank approximation in the field of machine learning. We think that the hidden neuron number should not be the latent dimension of the input, otherwise only local optima is reached and it is too greedy to be helpful to the whole network structure optimization. So we address the DNN network structure global optimization problem with stacked autoencoder structure, which is seldom investigated by other researchers with gradient descent-based learning algorithm. When stacked autoencoders are applied to construct a DNN, how many hidden layers should we need for a given data set? This DNN architecture design problem is seldom addressed in the literature, while in this work we propose the dynamical growth strategy to determine the

depth of the DNN architecture, which is quite different with pre-design network structure. Also, we use early stop regularization technique to find the optimal depth of the DNN, it can be explained as data driven architecture design.

Stacked autoencoders based DNN has the drawback for global optima. When we apply sparse or IB constraints to learn data representation, if the proceeding layer output has sparse /low rank, it will become more difficult for succeeding layers to learn more sparse /low rank representation again. In addition to the strategy to avoid local optima in the proposed learning scheme, other solution is to adopt stacked generalization technique [24] to construct deep and wide learning system for further improving the system performance on given task.

#### D.2 Discussion of Related Work

(a) **Related with other non-gradient descent algorithm:** Recently, it is found that under the name of extreme learning machine (ELM), several PIL variants related works have been published in this journal [37][38][39]. In fact, ELM is a renamed work, it originates from the PIL, detailed explanations can be found in [31] and more discussions will be presented in the later works.

(b) **Related with principal component analysis (PCA):** PCA is the optimal linear autoencoder [32][33][34][42]. SVD method can be used to calculate pseudoinverse or used in PCA. Bourlard *et al* believed that for autoencoder the nonlinearities of the hidden units are useless and that the optimal parameter values can be derived directly by SVD and low rank matrix approximation. Their approach appears only related linear case, it is a PCA approximation and encoder weight is the same as decoder weight [33]. However, in our proposed method, we apply truncated SVD (low rank approximation) to pseudoinverse of input matrix to compute encoder weight matrix, and then apply nonlinear transformation to obtain hidden layer output.

(c) **Related with our previous work:** In our previous work [24], we set hidden neuron number to be equal to the training sample number in order to realize exact learning, while in this work, low rank approximation is adopted for representation learning. This work is the extension of paper [25], here global optimization of the DNN structure is investigated. We find that Wong *et al*'s work [37] is more similar with our previous work [24] in learning strategy, the main difference is that activation function is taken as kernel function in their work. Detailed analysis about the similarity and difference will be presented in later work, here we just point out that random projection cannot learn representation from data.

#### D.3 Further work

Currently, the PILAE algorithm can be applied when the loss error obeys Gaussian density distribution, other probabilities such as Laplace probability should be studied in the further work.

For the desired output vector  $x$ , the probability model of an autoencoder can be expressed as follows:

$$x = W_d h + \varepsilon, \quad (38)$$

where  $h$  is the output vector of hidden neurons:  $h = f(W_e x)$ ,  $\varepsilon$  is noise. Suppose  $\varepsilon \sim N(0, \sigma_1^2 I)$ , with Eq. (38):

$$p(x | W_d h) = \frac{1}{(2\pi)^{d/2} \sigma_1^d} \exp\left(-\frac{1}{2\sigma_1^2} \|x - W_d h\|^2\right). \quad (39)$$

If the elements in the encoder matrix  $W_d$  are *i.i.d*, and the prior probability of the element weight  $w_j \sim \text{Laplace}(0, \sigma_2)$ ,

that is,  $p(w_j) = \frac{1}{2\sigma_2} \exp\left(-\frac{|w_j|}{\sigma_2}\right)$ , then the loss function  $L$  becomes:

$$L = \sum_{i=1}^N \|x_i - W_d h\|^2 + \frac{\sigma_1^2}{\sigma_2^2} \sum_{j=1}^{dp} |w_j| \quad (39)$$

This loss function is the sparse coding of an autoencoder, it is also called Lasso in statistics. The proposed PILAE algorithm will be developed for solving the above  $L_1$  norm constraint problem in the future.

Here we should note that hidden neuron number and hidden layer number are two hyperparameters for DNN architecture, it is still an open problem to realize AutoML. Currently, the methods to find an optimal DNN topology architecture include heuristic design, or search in hyper-parameter space. Our work belongs to the later one, but with empirical estimation formula for reducing the region of search space. This work is our effort toward to AutoML.

## V. CONCLUSION

This paper focuses on developing a non-gradient descent deep learning algorithm, which can reduce the training time for deep feedforward neural network. This is achieved by a feedforward learning algorithm based on pseudoinverse theory and low rank approximation. No iterative optimization is needed for training DNN, and the proposed pseudoinverse learning algorithm with low rank approximation can greatly speed up the learning processing. It is shown by experiments that the proposed learning scheme can reduce the training time drastically. The representation learning performance of the networks with the proposed learning scheme is better than the BP algorithm based network on the several datasets. The experiments illustrate that our proposed algorithm is a better scheme for training DNN on considering both the time consuming and representation learning capacity.

## REFERENCES

- [1] Y. Bengio, A. Courville, P. Vincent, "Representation learning: a review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798-1828, 2013.
- [2] Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton, "ImageNet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems* 25 (NIPS 2012), 2012.
- [3] Dan Ciresan, Ueli Meier, and Jürgen Schmidhuber, "Multi-column deep neural networks for image classification," in the proceedings of the 2012 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [4] M. Gong, J. Zhao, J. Liu, Q. Miao, and L. Jiao, "Change detection in synthetic aperture radar images based on deep neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 27, no. 1, pp. 125–138, Jan 2016.
- [5] W. Hou, X. Gao, D. Tao, and X. Li, "Blind image quality assessment via deep learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 26, no. 6, pp. 1275–1286, June 2015.
- [6] George E. Dahl, D. Yu, L. Deng, A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," *IEEE Transactions on Audio, Speech, and Language Processing*, 20 (1), 30–42, 2012.
- [7] H.-S. Le, I. Oparin, A. Allauzen, J.-L. Gauvin, and F. Yvon, "Structured output layer neural network language models for speech recognition," *IEEE Trans. Audio, Speech and Language Processing*, 21(1), 197-206, 2013.
- [8] Hanna Kamyshanska and Roland Memisevic, "The Potential Energy of an Autoencoder," *IEEE Transaction on pattern analysis and machine intelligence*, vol. 37, No. 6, 1261-1273, 2015.
- [9] Geoffrey E. Hinton and Ruslan R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science* 313.5786, 504-507, 2006.
- [10] Yoshua Bengio, "Learning Deep Architectures for AI," *Foundations and Trends in Machine Learning*, Vol. 2, No. 1, 1-127, 2009.
- [11] M. Ranzato, C. Poultney, S. Chopra, and Y. LeCun, "Efficient learning of sparse representations with an energy-based model," *Advances in Neural Information Processing Systems* 19, 1137-1144, 2007.
- [12] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Machine Learning Res.*, 11, 3371-3478, 2010.
- [13] S. Rifai, P. Vincent, X. Muller, X. Glorot, and Y. Bengio, "Contractive auto-encoders: Explicit invariance during feature extraction," In *ICML'2011*.
- [14] M. Denil, B. Shakibi, L. Dinh, & N. de Freitas, "Predicting parameters in deep learning," *Advances in Neural Information Processing Systems*, 2148-2156, 2013.
- [15] Dong Yu, Frank Seide, Gang Li, "Exploiting sparseness in deep neural networks for large vocabulary speech recognition," in the Proceedings of the 2012 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2012.
- [16] Shun-Ichi Amari, "Natural gradient works efficiently in learning," *Neural Comput.*, 10(2):251–276, 1998.
- [17] Kui Jia, "Improving training of deep neural networks via singular value bounding,," *CoRR*, abs/1611.06013, 2016 .
- [18] Tim Salimans and Diederik P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," In *NIPS*, 2016.
- [19] B. Chandra, Rajesh K. Sharma, "Fast learning in Deep Neural Networks," *Neurocomputing*, Vol. 2, No.171, 1205-1215, 2016.
- [20] Shai Shalev-Shwartz, Ohad Shamir, and Shaked Shammah, "Failures of Gradient-Based Deep Learning," arXiv:1703.07950, 2017.
- [21] P. Guo, C.L.P. Chen and Y.G. Sun, "An Exact Supervised Learning for a Three-Layer Supervised Neural Network", *Proceedings of the International Conference on neural Information Processing (ICONIP'95)*, pp1041-1044, Beijing, 1995.
- [22] P. Guo and M. Lyu, "Pseudoinverse Learning Algorithm for Feedforward Neural Networks," in *Advances in Neural Networks and Applications*, (N.E. Mastorakis, Ed.), Puerto De La Cruz, Tenerife, Canary Islands, Spain, Feb. , 321-326, 2001.
- [23] P. Guo, M.R. Lyu and C.L.P. Chen, "Regularization Parameter Estimation for Feedforward Neural Networks ", *IEEE trans System, Man and Cybernetics* (B), Vol.33, No.1, pp.35-44, 2003.
- [24] P. Guo and M. Lyu, "A pseudoinverse learning algorithm for feedforward neural networks with stacked generalization applications to software reliability growth data," *Neurocomputing*, vol. 56, 101-121, 2004.
- [25] K. Wang, P. Guo, X. Xin and Z. B. Ye, "Autoencoder, Low Rank Approximation and Pseudoinverse Learning Algorithm", *Proceedings of the 2017 IEEE International Conference on System, Man and Cybernetics*, pp.948--953, 2017.
- [26] Markovsky Ivan, "Low rank approximation: algorithms, implementation, applications," Springer Science & Business Media, 2011.
- [27] T. L. Boullion, P. L. Odell, "Generalized Inverse Matrices," Wiley, New York, 1971
- [28] G. H. Golub, C. F. Van Loan, "Matrix Computations," John Hopkins Uni. Press, 4th edition, 2013.
- [29] H. Narayanan and S. Mitter, "Sample complexity of testing the manifold hypothesis," In J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, 23, 1786–1794, 2010.
- [30] C.M. Bishop, "Neural Networks for Pattern Recognition," Cambridge, 2005.
- [31] P. Guo, "A Vest of the Pseudoinverse Learning Algorithm", Preprint, <http://dx.doi.org/10.13140/RG.2.2.34706.56005/2>, 2018.
- [32] P. Baldi and K. Hornik. "Neural networks and principal component analysis: Learning from examples without local minima". *Neural Networks*, 2:53–58, 1988.
- [33] H. Bourlard and Y. Kamp, Auto-Association by Multilayer Perceptrons and Singular Value Decomposition, *Biological Cybernetics*, 59, 291-294, 1988.
- [34] E. Oja. Principal components, minor components and linear neural networks. *Neural Networks*, 5:927–935, 1992.
- [35] Lok-Won Kim, "DeepX: Deep Learning Accelerator for Restricted Boltzmann Machine Artificial Neural Networks", *IEEE Transactions on Neural Networks and Learning Systems*, Vol. 29, Iss.5, pp. 1441 – 1453, 2018.

- [36] Christian P. Hansen, “The truncated SVD, as a method for regularization”. BIT Computer Science and Numerical Mathematics, 1987.
- [37] Chi Man Wong; Chi Man Vong; Pak Kin Wong; Jiuwen Cao, “Kernel-Based Multilayer Extreme Learning Machines for Representation Learning”, IEEE Transactions on Neural Networks and Learning Systems; Vol. 29, Iss.3, pp. 757 – 762, 2018.
- [38] Lei Zhang; David Zhang, “Evolutionary Cost-Sensitive Extreme Learning Machine”, IEEE Transactions on Neural Networks and Learning Systems, Vol.28, Iss. 12, pp. 3045 – 3060, 2017.
- [39] Mingxing Duan; Kenli Li; Xiangke Liao; Keqin Li, “A Parallel Multiclassification Algorithm for Big Data Using an Extreme Learning Machine”, IEEE Transactions on Neural Networks and Learning Systems; Volume 29, Issue 6, pp. 2337 - 2351, June 2018.
- [40] B. O. Ayinde, J. M. Zurada, “Deep Learning of Constrained Autoencoders for Enhanced Understanding of Data”, IEEE Transactions on Neural Networks & Learning Systems, Date of Publication: 26 September 2017.
- [41] G Alain, Y Bengio, S Rifai “Regularized auto-encoders estimate local statistics,” arXiv:1211.4246 [cs.LG], 2012.
- [42] M. Magdonismail, C. Boutsidis, “Optimal Sparse Linear Auto-Encoders and Sparse PCA”, Advances in Neural Information Processing Systems 29, pp. 298—306, 2016.
- [43] J. Vanschoren, J. N. van Rijn, B. Bischl, and L. Torgo. “OpenML: networked science in machine learning”. SIGKDD Explorations., vol. 15, no. 2, pp. 49–60, 2013.
- [44] M. Lichman, “UCI machine learning repository,” School Inf. Comput. Sci., Univ. California, Irvine, CA, USA, Tech. Rep., 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [45] K. Wang, P. Guo, Q. Yin, et al, “A pseudoinverse incremental algorithm for fast training deep neural networks with application to spectra pattern recognition”, *Proceedings of the 2016 International Joint Conference on Neural Networks*, (IJCNN 2016), pp.3453-3460, 2016.
- [46] X. Han, K. Rasul, and R. Vollgraf. “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms”, *arXiv preprint arXiv:1708.07747* (2017).