

Information Technology and Quantitative Management (ITQM 2017)

Hybrid Learning Network: A Novel Architecture for Fast Learning

Ying Liu^{a,b}, Chao Xiang^a

^a*School of Computer and Engineering, University of Chinese Academy of Sciences, Beijing, 100190 China*

^b*Key Lab of Big Data Mining and Knowledge Management, Chinese Academy of Sciences, Beijing, 100190 China*

Abstract

There're many effective architectures of the artificial neural network(ANN). For which the training is a hard work. The cost for training an ANN increases exponentially when the ANN gets deeper or wider. We therefore propose a novel architecture, the Hybrid Learning Network(HLN), to achieve a fast learning with good stability. The HLN can learn from both labeled data and unlabeled data at the same time in a hybrid learning manner. It uses a Self Organizing Map unified by the specially designed nonlinear function as the sparsity mask for a hidden layer to improve the training speed. We experiment our architecture on a synthetic dataset to test its regression capability against the traditional architecture, the result is promising. We also experiment on the well-known MNIST dataset, it demonstrates an even more impressive learning efficiency by up to 40× speed-up on training. It indicates that our architecture performs better on handling larger datasets and tackling larger problems.

© 2017 The Authors. Published by Elsevier B.V.

Selection and/or peer-review under responsibility of ITQM2017.

Keywords: HLN; hybrid learning; neural network; sparsity mask; SOM; MNIST; fast learning

1. Introduction

Since the first mathematical model of the artificial neural network(ANN) was proposed in 1943[1], lots of different architectures have been proposed to develop the model, such as the Convolutional Neural Network(CNN) for image recognition [2], the Recurrent Convolutional Neural Network(R-CNN) for object detection in videos [3], and the Long Short Term Memory(LSTM) for speech recognition [4], etc. These specially designed neural networks are trained by a few efficient methods such as the stochastic gradient descent(SGD)[5]. Such architectures suit well for their specific applications but may have plain or worse performances on others, and their best performances rely heavily on the hyperparameter configuration. Therefore a relatively universal architecture that enables equal or similar performances among varied applications is in real demand.

Although a number of different neural network architectures have been proposed in the past years, literally all of them can be classified into 3 categories, supervised [6], unsupervised [7], and semi-supervised learning [8]. In supervised learning, one can only train a model on labeled samples. However in real applications, labeling a large dataset is infeasible as a result unlabeled samples are the major. To make use of unlabeled data, a few unsupervised

*Corresponding author. Tel.: +86-183-0114-2368;

E-mail address: xiangchao215@mails.ucas.ac.cn.

training algorithms emerged. These unsupervised learning methods attempt to produce an optimized parameter initialization [9]. Thus, such unsupervised techniques can only be applied before the supervised training phase. Once the supervised begins, unsupervised learning will quit. Semi-supervised learning allows the model to learn from both labeled and unlabeled samples at the same time, where a regularizer is embedded in the supervised optimizing object, and generally a balance constraint is required to avoid the trivial solution [10]. Such an enhanced learner brings more hyperparameters (such as the balance constraint), making it even more difficult to search for the best settings for current architectures. Additionally, this integration manner makes it impossible to separate the two learners as the semi-supervised regularizer is built on the supervised learner, and therefore an early supervised training alone with profound labeled data is necessary before the semi-supervised learning starts.

Therefore in the paper, we introduce a new architecture that combines the supervised and unsupervised learning requiring no extra techniques in training. Our architecture, the Hybrid Learning Network (HLN), consists of stacked hybrid learning layers, each of which embeds a Self Organizing Map (SOM). In the similar manner of backpropagation, HLN demonstrates higher learning capability and stability.

The main contributions of our work are:

- HLN is proposed as a hybrid learning architecture to learn both unlabeled and labeled data simultaneously. HLN overcomes the problem of traditional semi-supervised learning methods which requires an early standalone training for the supervised learner.
- A SOM-embedding layer structure is designed to learn a cluster mapping function from unlabeled data to speed up the supervised learning from labeled data.
- A nonlinear function $h(x)$, is proposed to measure the state of a SOM in each iteration of the training, and then determine a sparsity mask for every hidden layer.

Our proposed HLN is implemented in Python and all our code and results of experiments are available at <https://github.com/hiroki-kyoto/hybrid-learning-net>.

The rest of the paper is organized as follows. In section 2 we introduce existing semi-supervised algorithms for neural network models and SOM. In section 3, we present our novel architecture HLN and how to embed SOMs into deep architectures of neural networks. In section 4 we explain the training theory for HLN. Section 5 presents the experimental results, and the last section concludes our work.

2. Related work and background Knowledge

2.1. Semi-supervised learning for neural network

A key assumption in semi-supervised algorithms, is the structure assumption: two samples with similar distribution on the same mapping structure tend to have high probability belonging to the same class. Based on this assumption, one can use large unlabeled data to uncover such structures. There're already a few algorithms dedicated to it, such as cluster kernels [11], Low Density Separation (LDS) [12], label propagation [13], etc. In such algorithms, designing a regularizer to enable the model to learn the representation or structure of raw data becomes the key point.

Let's firstly focus on the general algorithm description of semi-supervised learning. Given a set of unlabeled samples, $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} (\mathbf{x} \in \mathbb{R}^d)$, and the similarity labels between any \mathbf{x}_i and \mathbf{x}_j , $\mathbf{W} = \{W_{ij} | i, j = 1, \dots, N\}$, our goal is to find the best embedding function, $f(\mathbf{x})$, for each sample \mathbf{x}_i , in order to minimize:

$$\Delta_f = \sum_{i=1}^N \sum_{j=1}^N L(f(\mathbf{x}_i), f(\mathbf{x}_j), W_{ij}) \quad (1)$$

where,

- $L(\cdot)$ is the loss function of 3 variables: $\langle f(\mathbf{x}_i), f(\mathbf{x}_j), W_{ij} \rangle$. For example, if Euclidean distance is used,

$$L(f(\mathbf{x}_i), f(\mathbf{x}_j), W_{ij}) = (\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\| - W_{ij})^2$$

- $f(\mathbf{x}) \in \mathbb{R}^n$ is the embedding function, trying to produce an output vector for \mathbf{x}_i , similar to that for \mathbf{x}_j with $W_{ij} = 0$, and dissimilar with $W_{ij} = 1$.
- $W_{ij} \in \mathbb{R}$ is the similarity label of the sample pair $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ from \mathbf{X} .

As a result, the overall optimization object is

$$\arg \min \sum_i \ell(f(\mathbf{x}_i), \hat{\mathbf{y}}_i) + \lambda \Delta_f \quad (2)$$

where $\ell(\cdot)$ is the loss function for supervised learning. λ is proposed as the balance coefficient to embed the semi-supervised regularizer Δ_f .

2.2. Self Organizing Map

Self Organizing Map is an effective method for the visualization of high-dimensional data. It can also be used as an automatic clustering method. The SOM consists of a two-dimensional regular grid of nodes. The models are automatically organized into a meaningful two-dimensional order in which similar models are closer to each other in the grid than the more dissimilar ones[14]. Rules used to update models are:

$$\mathbf{m}_i(t+1) = \mathbf{m}_i(t) + h_{c(x),i}(\mathbf{x}(t) - \mathbf{m}_i(t)) \quad (3)$$

and the learning rate is dynamically determined as

$$h_{c(x),i} = \alpha(t) \exp\left(-\frac{\|\mathbf{r}_i - \mathbf{r}_c\|^2}{2\sigma^2(t)}\right) \quad (4)$$

where $\mathbf{m}_i \in \mathbb{R}^n$ denotes the i^{th} model vector, \mathbf{x} is an input pattern, $c(x)$ relates to the best match vector index in \mathbf{m} for input pattern \mathbf{x} , and $\alpha(t)$ is a learning rate that decreases as the training proceeds. \mathbf{r} denotes the model vector location in the map, and $\sigma(t)$ corresponds to the width of the neighborhood function, which also decreases monotonically with the regression steps.

3. Hybrid Learning Network

Hybrid Learning Network is a novel architecture aiming at enhancing arbitrary networks on their training efficiency and stability. Each hidden layer in HLN architecture is embedded with a SOM. For the simplest case where neurons are fully connected, the traditional architecture, called Fully Connected Network(FCN) architecture, is presented in Figure 1(a). The proposed HLN architecture is presented in Figure 1(b).

In HLN architecture, we propose $h(x)$, an unifying function for SOM, to convert pattern dissimilarity $\|\mathbf{x} - \mathbf{m}_i\|$ to a semi-supervised learning factor for different hidden units. The semi-supervised learning factors act as a dynamic neuron activation sparsity mask for each hybrid learning layer. It works in the similar manner as the Dropout technique[15], enabling the neural networks to improve the model robustness and prevent overfitting by learning their submodels for each batch. However, HLN differs from techniques like Dropout because HLN does not generate sparsity with randomness. HLN uses the SOM to unsupervisedly learn the *static* policy of the neuron-activation distribution for each hidden layer, the network sparsity generator thereby will stabilize with training steps, and the randomness in sparsity will disappear automatically. We propose $h(x)$ as follows,

$$\delta_{max} = \max_i (\|\mathbf{m}_i - \mathbf{x}\| + \epsilon) \quad (5)$$

$$\delta_{min} = \min_i (\|\mathbf{m}_i - \mathbf{x}\| + \epsilon) \quad (6)$$

where δ_{max} is the maximum dissimilarity between an arbitrary vector \mathbf{m}_i and the input pattern vector \mathbf{x} , δ_{min} is the minimum one. The scale of all dissimilarities is

$$\delta_{scale} = \frac{\delta_{min}}{\delta_{max} - \delta_{min} + \epsilon} \quad (7)$$

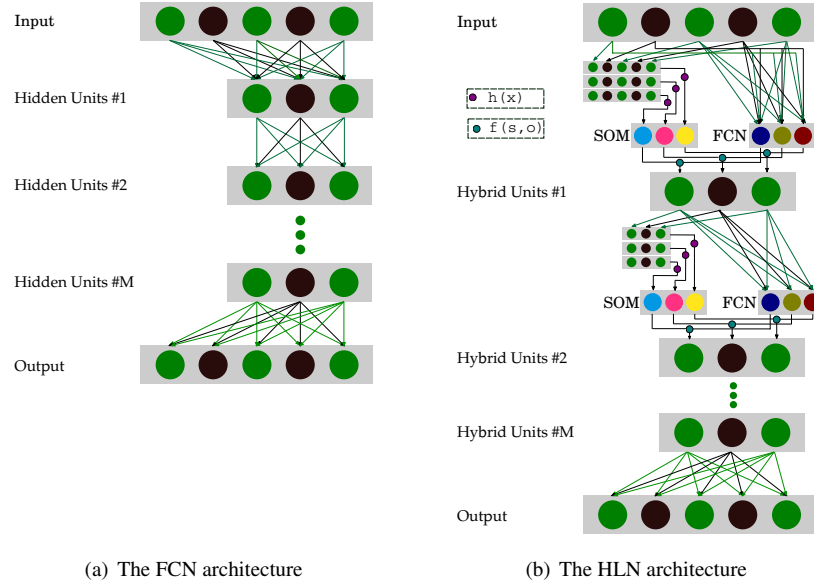


Fig. 1. The left is FCN architecture, the right is HLN architecture. The big solid circles are the neurons, the medium ones are the middle states of neurons, and the small ones are the vector maps \mathbf{m}^k of SOMs, finally the tiny ones are the function modules. The parts omitted refers to repeating hidden layers connected in same way.

We then unify the dissimilarities and convert them into sparsity mask as

$$h(x)|_{x=\|\mathbf{m}_i-\mathbf{x}\|} = \frac{\delta_{max} \cdot \delta_{scale}}{\|\mathbf{m}_i - \mathbf{x}\| + \epsilon} - \delta_{scale} \quad (8)$$

where \mathbf{m}_i and \mathbf{x} is the same way defined as in equation (3). ϵ is set at a small constant value, such as $\epsilon = 10^{-5}$, to avoid *division-by-zero* errors, which rarely happens. Notice that if any other metric is preferred, we can always replace the pattern dissimilarity $\|\mathbf{m}_i - \mathbf{x}\|$ with other forms. $f(s, o)$ is the function to combine the sparsity mask s generated by $h(\cdot)$ and the fully connected linear summation output o . In most cases, we use the simple multiplication as the operator,

$$f(s, o) = s \cdot o = h(\|\mathbf{m}_i - \mathbf{y}^{k-1}\|) \cdot \sum_j (w_j^{k,i} y_j^{k-1}) + b^{k,i}, k > 1 \quad (9)$$

As shown in Figure 1(b), $f(\cdot)$ is followed by an activation function, which could be any arbitrary nonlinear function that takes only one dimension input and outputs a single real value, such as the sigmoid function, $\tanh(x)$ and ReLU. With the new architecture, we update the forward computing rules as

$$y_i^k(\mathbf{x}) = \sigma \left(h(\|\mathbf{m}_i^k - \mathbf{y}^{k-1}(\mathbf{x})\|) \cdot \sum_j (w_j^{k,i} y_j^{k-1}(\mathbf{x})) + b^{k,i} \right), k > 1 \quad (10)$$

and for the first hybrid learning layer,

$$y_i^1(\mathbf{x}) = \sigma \left(h(\|\mathbf{m}_i^1 - \mathbf{x}\|) \cdot \left(\sum_j w_j^{1,i} x_j \right) + b^{1,i} \right) \quad (11)$$

where m_i^k denotes the i^{th} vector in the k^{th} SOM of the net(all indexes start from 1), and σ is a nonlinear activation function.

4. HLN Training

Let's compare the theories between HLN and the general model mentioned in section 2. Existing semi-supervised algorithms include a regularizer into the supervised learning optimizing object, making it impossible to train both of the supervised learner and the unsupervised one separately. As analyzed before, performance of such regularizer solutions depends on the standalone supervised pretraining which requires profound labeled data. However, in our proposed architecture HLN, we assign each of the learning methods a completely separated optimizing object, with no priority orders restricted.

For the supervised learning, we may have such an optimizing object,

$$\arg \min_g \sum_{i=1}^L \ell(g(\mathbf{x}_i), \hat{\mathbf{y}}_i) \quad (12)$$

where $g(\mathbf{x})$ is the function describing the mapping from the input $\mathbf{x} \in \mathbb{R}^d$ to the output $\mathbf{y} \in \mathbb{R}^n$, parameterized with the neural connection weights \mathbf{W} and the biases \mathbf{b} for all non-input layer neurons in the HLN architecture. Equation 12 becomes the Equation 13 when Enclidean metric is applied as the loss function.

$$\arg \min_{\mathbf{W}, \mathbf{b}} \frac{1}{2} \sum_{i=1}^L \|y^k(\mathbf{x}_i) - \hat{\mathbf{y}}_i\|^2 \quad (13)$$

For the unsupervised learning, the optimizing objects are

$$\arg \min_{\mathbf{m}^k} \sum_i^{L+U} \left(\min_j (\|\mathbf{m}_j^k - \mathbf{x}_i^k\|) \right), \quad k = 1, 2, \dots \quad (14)$$

with the predefined notation:

$$\mathbf{x}_i^k = \begin{cases} \mathbf{x}_i, & k = 1 \\ y^{k-1}(\mathbf{x}_i), & k > 1 \end{cases}, \quad \mathbf{x}_i \in L + U \quad (15)$$

For a net with M hybrid learning layers, there should be 1 supervised optimizing object and M unsupervised optimizing objects. Although a multi-objective problem may not have a global solution, in HLN, since each object is optimized in its own isolated parameter space, each optimizing object has its own global optimization solution.

Using the well-known Back Propagation(BP) algorithm, we can obtain parameter updating rules for HLN training. We focus on the difference between HLN and the original BP algorithm for the FCN model. Let's denote the linear summation for i^{th} unit in k^{th} layer as o_i^k , then the linear summations of hybrid learning layers are (with $\mathbf{y}^0 = \mathbf{x}$):

$$o_i^k = h(\|\mathbf{m}_i^k - \mathbf{y}^{k-1}\|) \sum_j W_j^{k,i} y_j^{k-1}, \quad k = 1, 2, \dots, M \quad (16)$$

Apply the non-linear activation function, for example, the sigmoid function:

$$\frac{\partial E}{\partial o_i^k} = \frac{\partial E}{\partial y_i^k} \sigma'(o_i^k) = \frac{\partial E}{\partial y_i^k} y_i^k (1 - y_i^k), \quad k = 1, 2, \dots, M + 1 \quad (17)$$

Thus the recursive layer gradient computing is shown in Equation 18.

$$\frac{\partial E}{\partial y_i^{k-1}} = \sum_j \frac{\partial E}{\partial o_j^k} h(\|\mathbf{m}_j^k - \mathbf{y}^{k-1}\|) W_j^{k,i} \quad k = 1, 2, \dots, M \quad (18)$$

The partial error gradients for the connection parameters \mathbf{W} are shown in Equation 19.

$$\frac{\partial E}{\partial W_j^{k,i}} = \sum_i \frac{\partial E}{\partial o_i^k} h(\|\mathbf{m}_i^k - \mathbf{y}^{k-1}\|) y_i^{k-1}, \quad k = 1, 2, \dots, M \quad (19)$$

The partial error gradients on biases are presented in Equation 20.

$$\frac{\partial E}{\partial b^{k,i}} = \frac{\partial E}{\partial o_i^k} \frac{\partial o_i^k}{\partial b^{k,i}} = \frac{\partial E}{\partial o_i^k} \times 1, \quad k = 1, 2, \dots, M+1 \quad (20)$$

Therefore we can iterate recursively along layers using Equation (16)(17) (18)(19) (20) to update all the parameters of the supervised training, and the corresponding parameter updating rules for the unsupervised training are shown in Equations (3)(4).

5. Empirical Study

5.1. Synthetic data

A synthetic dataset is generated. It maps a set of 2-dimension vectors into 3 classes, where 3.5% of the noise is included to simulate a real application. We then apply the data to an instance of HLN and FCN. In both networks, we only use a single hidden layer. Figure 2-3 present the progress of the accuracy as the number of epoches increases.

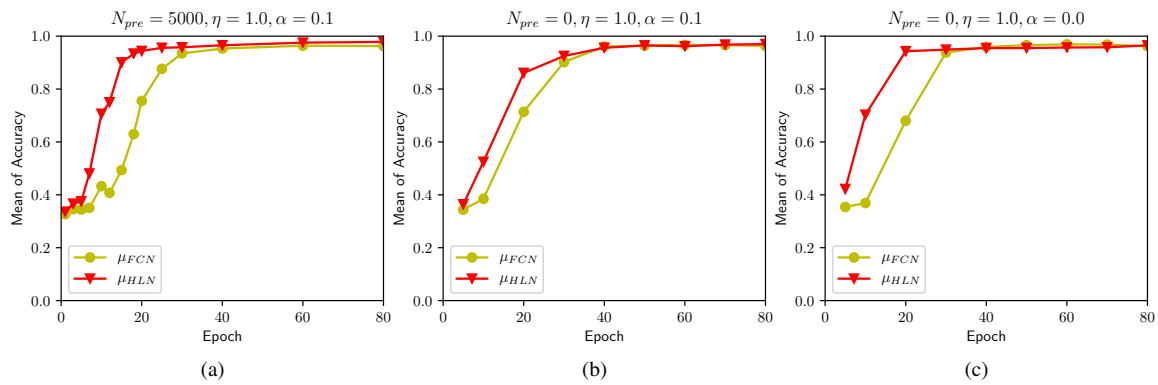


Fig. 2. Mean of accuracy on the synthetic data: μ is the mean accuracy for multiple replays with the same configuration. N_{pre} is the number of unlabeled samples for pretraining, η is the initial learning rate, α is the decay factor of learning rate $\eta(t)$.

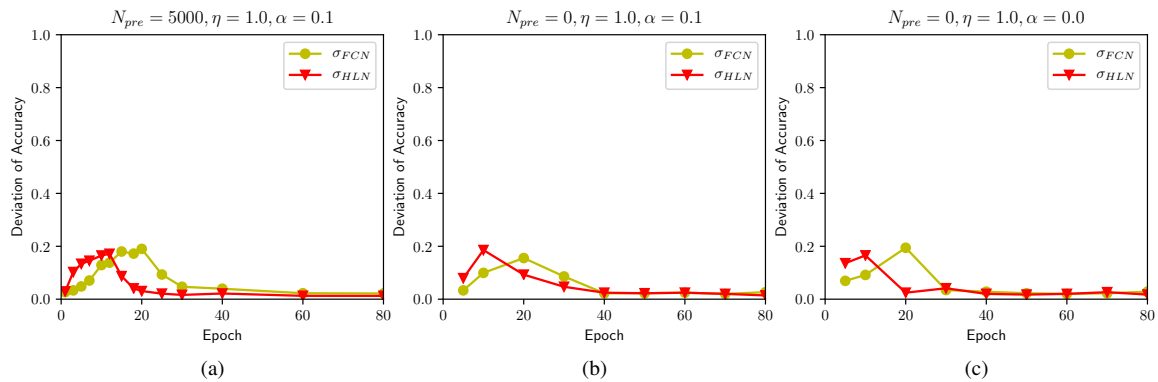


Fig. 3. Deviation of accuracy on the synthetic data: the deviation takes the form as δ is the deviation of multiple replays with the same configuration $\delta = \sum_1^N (\gamma_i - \mu)/N$, where γ_i is the accuracy of i^{th} replay.

In the experiment, we split the data into two partitions, one for the training of two architectures and the another for testing with the train/test ratio 0.8. Considering the possible impact from the parameter initialization,

we replay the complete workflow of training and testing with varied random parameter initializations for every net with different configurations. In Figure 2, It is evident that networks with our HLN architecture win over networks without HLN in all cases. The experimental results indicate:

- HLN architecture learns much faster than FCN.
- The pretraining with unlabeled data does improve the performance of HLN, though the improvement in accuracy is not evident.
- A constant learning rate is more advisable than a varying learning rate.

Figure 3 shows that HLN stablized quickly in all cases, while FCN requires more training to stablize. It indicates that HLN is more robust to parameter initialization.

5.2. MNIST

Our experiments on the well-known MNIST benchmark demonstrated even greater enhancement over the traditional architecture. Moreover the advantage of HLN is more evident with higher input dimension and deeper network. To enable the input dimension variation, a max-pooling is applied before the input of a neural net. The MNIST dataset provides 50,000 images for training and 10,000 images for testing, with each image resized as 28×28 of pixels. A max-pooling with the kernel of 2×2 leads to a input dimension of 196, and particularly a max-pooling with the kernel 1×1 will keep the original input dimension as 784. All networks used are configured as in Figure 1.

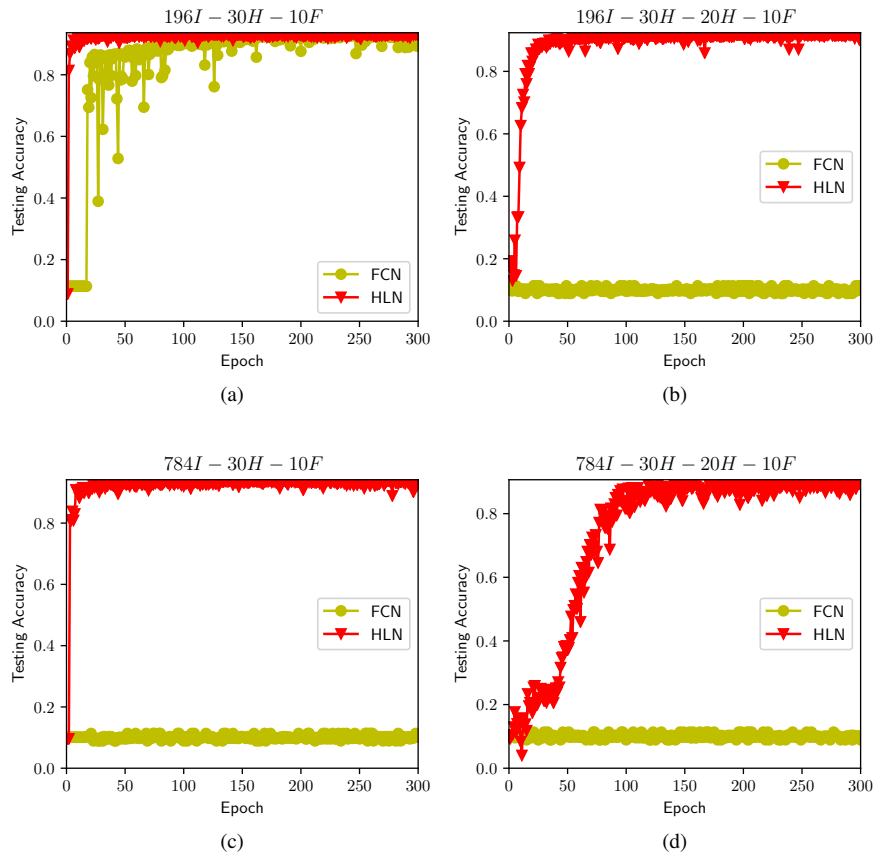


Fig. 4. Accuracy on MNIST with 4 different network models: 196I denotes the input layer of dimension 196, 30H denotes the hidden layer of dimension of 30, and 10F denotes the fully connected output layer of dimension 10.

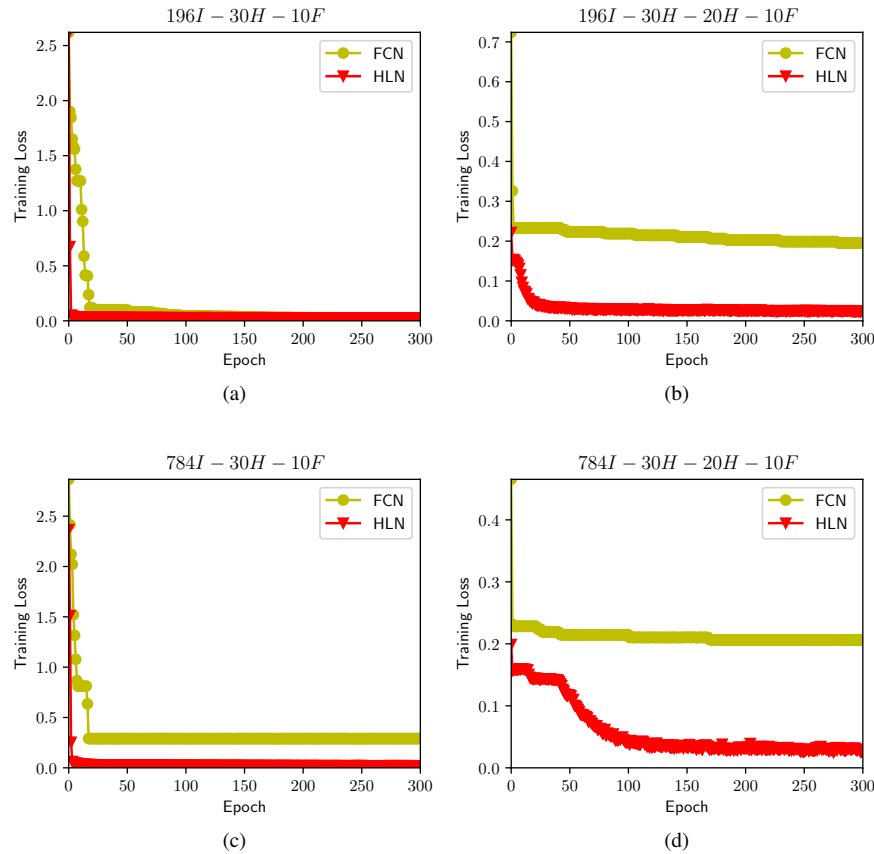


Fig. 5. Training loss on MNIST with 4 different network models

Observed from Figure 4(a), where a single hidden layer architecture is used, the FCN requires about 150-200 epoches to stabilize at its best test accuracy of 93%, however HLN only takes 4-5 epoches to achieve the same test accuracy or even higher. Thus it is a 40× speed up on the training without a sacrifice of the model capability. The HLN additionally provides much higher stability at all time during the entire training process. The reason behind is that HLN enables the network to unsupervisedly learn the neuron activation clustering map(the SOM) of every non-output layer especially the input layer, along with the supervised learning. Each sample for the input is more thoroughly taken in use, making HLN learn much faster.

In Figure 4(b) with a deeper architecture, FCN seems fail to learn at all. In contrast, HLN keeps a fast learning. Though HLN seems *slow down* a little bit for the sake of doubling the hidden layers. The observation follows: a deeper architecture brings a tougher parameter tuning.

In Figure 4(c), we use a max-pooling with a kernel of 1×1 , which equals to the original input as in MNIST. The input dimension of the network is therefore 4 times larger than that in Figure (4(a)). However, HLN learns as fast as the one with small input dimension with equal stability. Again FCN fails to learn.

In Figure 4(d), HLN takes about 120 epoches to stabilize at its best performance but FCN seems never start to learn. Compared to Figure 4(c), the networks have a deeper architecture, making the model more difficult to learn as expected. Compared to Figure 4(b), it has a larger input dimension, which leads to a tougher learning too. Putting two factors together turns out more harmful to the training as we may expect in a linear manner.

Figure 5 further confirms the facts observed from Figure 4. The test accuracy is not improved at all with deeper architectures, whereas the training loss is dropping, though, at a very slow pace. Figure 6 shows the neuron activation sparsity for every hidden layer in networks. Notice that a SOM-embedding output in HLN is used

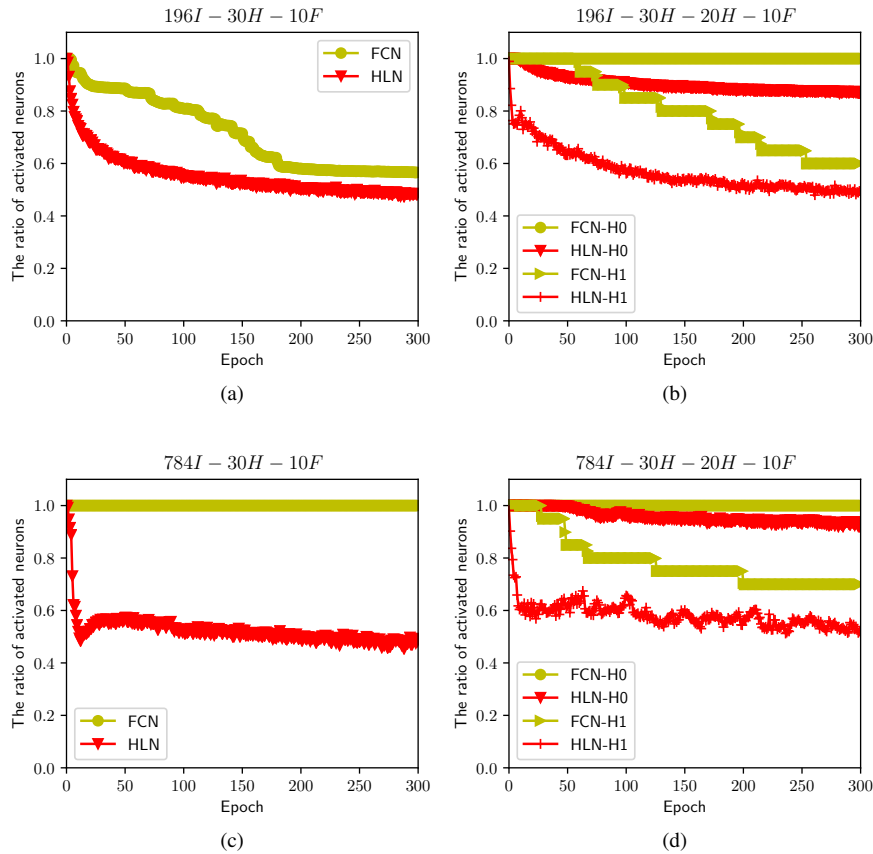


Fig. 6. Neuron activation sparsity on MNIST with 4 different network models: H0 denotes the first hidden layer, H1 denotes the second hidden layer.

as a sparsity mask to control, or more precisely, to increase the sparsity of every hidden layer. It's obvious that networks with HLN architecture tend to have a higher sparsity than the ones without.

6. Conclusions

We proposed HLN, a hybrid learning architecture for ANNs. HLN enables ANNs to learn from labeled and unlabeled data simultaneously and separately. We examined experiments on our synthetic dataset and the MNIST benchmark dataset. As the experiments showed, ANNs of HLN architecture proved better performance in training efficiency and stability. In the future, we will investigate other capabilities of HLN on more public datasets such as CIFAR-10 and CIFAR-100.

Acknowledgements

This project was partially supported by Grants from Natural Science Foundation of China #71671178/#91546201/#61202321, and the open project of the Key Lab of Big Data Mining and Knowledge Management. It was also supported by Hainan Provincial Department of Science and Technology under Grant No. ZDKJ2016021, and by Guangdong Provincial Science and Technology Project 20162016B010127004.

References

- [1] W. S. McCulloch, W. Pitts, A logical calculus of the ideas immanent in nervous activity, *Bulletin of mathematical biology* 5 (4) (1943) 115–133.
- [2] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [3] R. Girshick, Fast r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [4] A. Graves, N. Jaitly, A.-r. Mohamed, Hybrid speech recognition with deep bidirectional lstm, in: *Automatic Speech Recognition and Understanding (ASRU)*, 2013 IEEE Workshop on, IEEE, 2013, pp. 273–278.
- [5] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: *Proceedings of COMPSTAT'2010*, Springer, 2010, pp. 177–186.
- [6] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, L. D. Jackel, Handwritten digit recognition with a back-propagation network, in: *Advances in neural information processing systems*, 1990, pp. 396–404.
- [7] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion, *Journal of Machine Learning Research* 11 (Dec) (2010) 3371–3408.
- [8] O. Chapelle, B. Scholkopf, A. Zien, Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews], *IEEE Transactions on Neural Networks* 20 (3) (2009) 542–542.
- [9] Q. V. Le, Building high-level features using large scale unsupervised learning, in: *Acoustics, Speech and Signal Processing (ICASSP)*, 2013 IEEE International Conference on, IEEE, 2013, pp. 8595–8598.
- [10] R. Socher, J. Pennington, E. H. Huang, A. Y. Ng, C. D. Manning, Semi-supervised recursive autoencoders for predicting sentiment distributions, in: *Proceedings of the conference on empirical methods in natural language processing*, Association for Computational Linguistics, 2011, pp. 151–161.
- [11] O. Chapelle, J. Weston, B. Schölkopf, Cluster kernels for semi-supervised learning, in: *Advances in neural information processing systems*, 2003, pp. 601–608.
- [12] O. Chapelle, A. Zien, Semi-supervised classification by low density separation., in: *AISTATS*, 2005, pp. 57–64.
- [13] X. Zhu, Z. Ghahramani, Learning from labeled and unlabeled data with label propagation.
- [14] T. Kohonen, The self-organizing map, *Neurocomputing* 21 (1) (1998) 1–6.
- [15] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: a simple way to prevent neural networks from overfitting., *Journal of Machine Learning Research* 15 (1) (2014) 1929–1958.