

# Pseudo Boosted Deep Belief Network

Tiehang Duan and Sargur N. Srihari

Department of Computer Science and Engineering  
The State University of New York at Buffalo  
Buffalo, NY 14260, United States  
tiehangd@buffalo.edu, srihari@cedar.buffalo.edu

**Abstract.** A computationally efficient method to improve classification performance of a Deep Belief Network(DBN) is introduced. In the Pseudo Boost Deep Belief Network(PB-DBN), top layers are boosted while lower layers of the base classifiers share weights for feature extraction. PB-DBN maintains the same time complexity as a DBN with fast convergence to optimality by introducing the mechanism of pseudo boost. Experiments in classification show that after only a few iterations, the PB-DBN has higher accuracy than a classic DBN.

## 1 Introduction

A Deep Belief Network(DBN) is a layered network formed with undirected connections between its top two layers and downward directed connections between all its lower layers (Fig. 1(a)) [13]. They are used to infer posterior probabilities and make generative predictions. Restricted Boltzmann Machines(RBM) formed by adjacent layers in DBN are pre-trained from bottom-up layer by layer, then the whole network is fine tuned with back propagation to serve discriminative purposes or up down algorithms for generative purposes [7]. Many approaches to improve the performance and time efficiency of DBNs have been published since G. Hinton et al. proposed the method in 2006 [7]. The improvements include new training methods such as re-weighted wake sleep algorithm [4], exponential loss function gradient method [11]; new architectures of the network such as convolutional architectures [9], sparse architectures [8], hierarchical architectures [15], sequential architectures [1]; and there are also works combining the power of other learning methods [16].

Boosting is an ensemble learning method that is used to improve the performance of base classifiers. Ping Liu et al. combined DBN with boosting for facial expression recognition [10], and revealed the potential to integrate boosting mechanism into DBN (Fig. 1(b)). In this paper, we integrate the power of boosting into DBN with the boosting mechanism re-weighting each data vector for each iteration and DBN re-represent the data on each layer. As the time complexity of boosting is polynomial to weak classifier's time complexity(the number of classifiers determine the order of polynomial), it is impractical to directly apply DBN to be the base classifier— since a single DBN is already computationally expensive. We address this issue with a computationally efficient

weight sharing boost mechanism, with the resulting time complexity being of the same order as that of a single DBN.

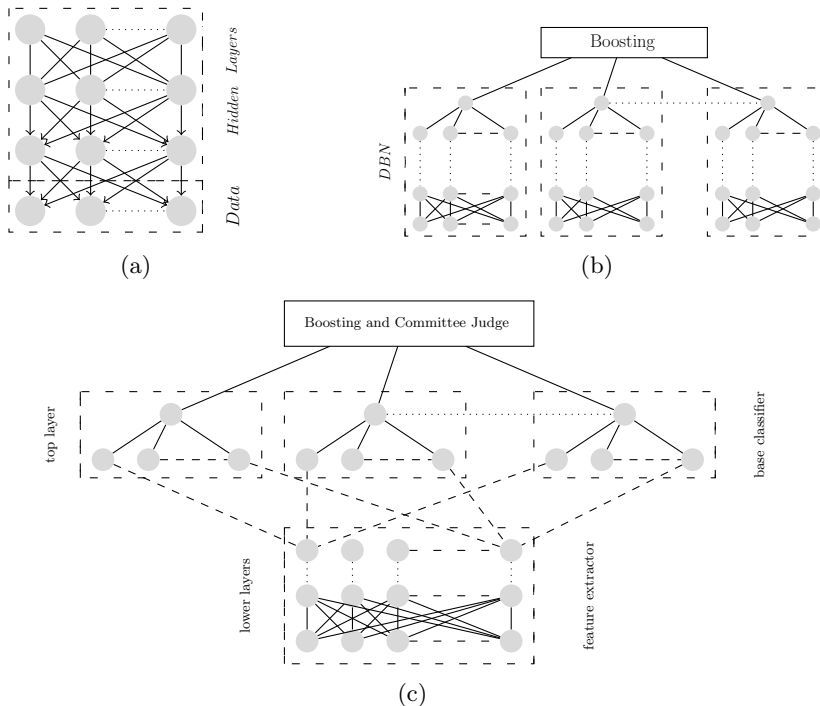


Fig. 1: Structure of (a) Deep Belief Network(DBN), (b) Boosted Deep Belief Network (BDBN) and (c) Pseudo Boosted Deep Belief Network(PB-DBN).

## 2 Pseudo Boosted Deep Belief Network

### 2.1 Initialization of the Network

Initialization of the proposed Pseudo-Boosted Deep Belief Network (PB-DBN) (Fig. 1 (c)) is analogous to that of a DBN. RBMs formed by adjacent layers are trained in turn from bottom up by minimizing the *contrastive divergence*:

$$KL(P^0 || P_{\theta}^{\infty}) - KL(P_{\theta}^n || P_{\theta}^{\infty}) \quad (1)$$

The first term in (1) is the Kullback Leibler divergence between the distribution of the data  $P^0$  and the equilibrium distribution of the model  $P_{\theta}^{\infty}$ . The second term is the KL divergence between the posterior distribution  $P_{\theta}^n$  and  $P_{\theta}^{\infty}$  [2]. To make the probabilistic model computationally tractable and also easy to fine

tune, we assume direct probabilistic dependence only exists between the adjacent layers, and the probability of activation for each hidden layer is the generalized linear sigmoid model of the adjacent hidden layer:

$$P(s_j = 1) = \frac{1}{1 + \exp(-b_j - \sum_i s_i w_{ij})} \quad (2)$$

where  $s_i$  is the node in adjacent layer,  $s_j$  is the node in current layer,  $w_{ij}$  is the corresponding edge weight and  $b_j$  is the bias term. The adjustment of weights during the pre-training phase follows the gradient of log probability of training data:

$$\frac{\partial \log(p(v^0))}{\partial w_{ij}} = \langle v_i^0 h_j^0 \rangle - \langle v_i^\infty h_j^\infty \rangle \quad (3)$$

where  $v_i^0$  is the node in the visible layer under initialization,  $v_i^\infty$  is the node in the visible layer after reaching equilibrium,  $h_j^0$  is the node in the hidden layer after initialization,  $h_j^\infty$  is the node in the hidden layer after reaching equilibrium,  $\langle v_i^0 h_j^0 \rangle$  is the correlation between  $v_i^0$  and  $h_j^0$ , and  $\langle v_i^\infty h_j^\infty \rangle$  is the correlation between  $v_i^\infty$  and  $h_j^\infty$  [5]. After the training of each RBM, the value in hidden layer is stored and used as data for the next layer's training.

When all RBMs finish training, the change of weights in the higher layers make the lower layers no longer optimal for the whole network [2]. Two fine-tuning methods can be used to further adjust weights. Error back propagation can be used for discriminative tuning and up down algorithms can be used for generative purposes. A detailed description can be found in [7].

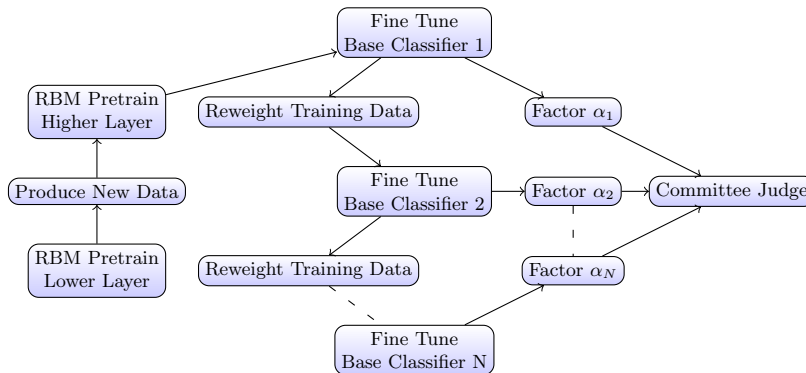


Fig. 2: Illustration of PB-DBN Mechanism.

## 2.2 Mechanism of Pseudo Boost

The boosting mechanism assigns different weights to data points based on the importance of each data [3]. The data weights are iteratively adjusted based on

classification result in previous iteration:

$$w_n^{(m+1)} = w_n^{(m)} \exp\left(\alpha_m I(y_m(\mathbf{x}_n) \neq t_n)\right) \quad (4)$$

where  $\mathbf{x}_n$  is the  $n$  th data,  $y_m(\mathbf{x}_n)$  is the classification decision made by the  $m$  th base classifier,  $t_n$  is the label of training data  $\mathbf{x}_n$ ,  $w_n^{(m)}$  is the weight of  $n$  th data in the  $m$  th iteration, and  $\alpha_m$  is defined as the log odd of the weighted error rate  $\varepsilon_m$ :

$$\varepsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (5)$$

$$\alpha_m = \frac{1}{2} \log\left(\frac{1 - \varepsilon_m}{\varepsilon_m}\right) \quad (6)$$

Then the updated data is used to train the next base classifier so the next base classifier can learn something new [14] [6]. And when the training phase is over, the final classification is made based on decision of the committee formed by the base classifiers:

$$Y_M(\mathbf{x}_n) = \text{sign}\left(\sum_{m=1}^M \alpha_m y_m(\mathbf{x}_n)\right) \quad (7)$$

where  $y_m(\mathbf{x}_n)$  is as defined in (4) and  $\alpha_m$  as defined in (6).

The model can be easily implemented with simple base classifiers such as Stumps, SVMs or simple decision trees [12]. However, it is met with computational complexity issues when a DBN is the base classifier. Since time complexity of boosting is polynomial to time complexity of each base classifier with the order of the polynomial being the number of classifiers. For most instances, it is impractical to directly adopt the classic boosting mechanism. Pseudo boosting can be viewed as a partial boost that operates only on the top layers of the network (Fig. 1 (c)).

In the proposed method, lower layers of the individual deep belief network serve as feature extractors and as they share great similarity to one another [7], it would be a waste of computational resource if we repeatedly train the lower layers every time we initialize a new base classifier, so every time when we start training the new classifier in PB-DBN, the lower layers of previous base classifier is directly used. In the fine tuning phase, different classifiers get different emphasis on the difficult training data(which is done by changing the data sampling distribution). Then the mechanism of boosting allows these classifiers to be unified together to minimize the total exponential error function:

$$E = \sum_{n=1}^N \exp\left(-t_n y_m(\mathbf{x}_n)\right) \quad (8)$$

where  $y_m(\mathbf{x}_n)$  and  $t_n$  are the same as defined in (4). Detailed illustration of the PB-DBN mechanism is shown in Fig. 2.

### 2.3 Time Complexity Analysis

We denote the training time needed for each RBM as  $t_{rbm}$ , the fine tuning time for the whole network as  $t_{tune}$ , and the time for classifier committee to make the judgment as  $t_c$ , then the training time of the PB-DBN model is:

$$T_{PB-DBN} = t_{rbm} + k \times t_{tune} + t_c \quad (9)$$

where  $k$  is the number of classifiers in the model. And if we adopt the normal boosting procedure, the training time would be:

$$T_{Boost-DBN} = k \times t_{rbm} + k \times t_{tune} + t_c \quad (10)$$

Comparing (9) and (10), the PB-DBN model reduces time complexity by avoiding repetitive training in the RBMs.

Table 1: Experiment Result on Different Network Settings

Hidden Layers	Iterations	Train Error Rate		Test Error Rate	
		DBN	PB-DBN	DBN	PB-DBN
(Lower -> Upper)	(Pretrain - Fine Tune)				
784 -> 256	6 - 6	2.58%	0.42%	3.66%	1.83%
784 -> 256	10 - 10	0.62%	0.12%	2.10%	1.76%
784 -> 256	20 - 20	0.02%	0.02%	1.78%	1.70%
784 -> 256	50 - 50	0.00%	0.00%	1.68%	1.60%
784 -> 784 -> 256	6 - 6	3.44%	0.36%	3.99%	1.86%
784 -> 784 -> 256	10 - 10	0.91%	0.20%	2.40%	1.72%
784 -> 784 -> 256	20 - 20	0.09%	0.07%	1.82%	1.65%
784 -> 784 -> 256	50 - 50	0.00%	0.00%	1.71%	1.64%

## 3 Experimental Results

As the performance of DBN is superior to boosting with normal base classifiers such as stumps, decision trees or single-hidden layer neural networks, here we compare the performance of PB-DBN with DBN.

We performed batch training with each batch containing 100 data and a total of 600 batches. Conjugate gradient descent is used in the fine tuning phase and we test our method with base classifier DBN of both 2 hidden layers and 3 hidden layers (Table 1). We see PB-DBN outperforms DBN on all sets of comparative experiments. The unique strength of PB-DBN method as revealed in Fig. 3 is its fast convergence to optimality, as the boosting mechanism is making the difficult data to be more significant and makes the training phase more efficient. The parameters including learning rate(0.1), momentum(0.2), weight initialization(0.1) and weight-decay(0.0002) are chosen by training the network several times and observing its performance on the separate validation dataset. Fig. 5 shows most

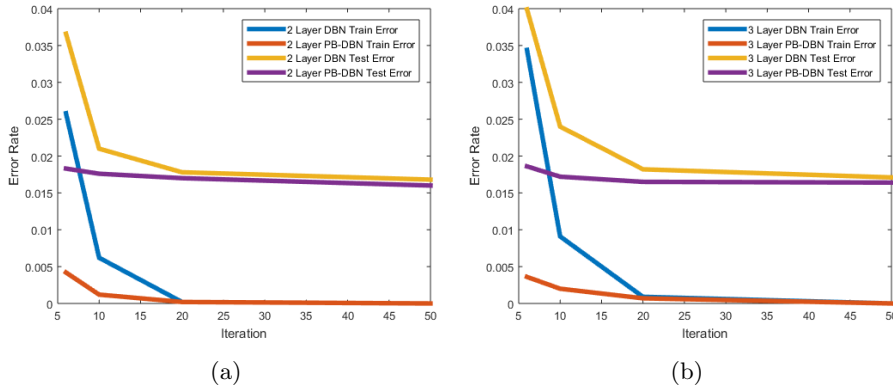


Fig. 3: Comparison of performance between DBN and PB-DBN. (a) Train error and test error of 2 layer networks. (b) Train error and test error of 3 layer networks.

Table 2: Performance with Different Number of Classifiers (10 Iterations)

Classifier Number	Hidden Layers	Train Error Rate	Test Error Rate
1	784 -> 784 -> 256	0.91%	2.40%
2	784 -> 784 -> 256	0.30%	1.96%
4	784 -> 784 -> 256	0.24%	1.88%
6	784 -> 784 -> 256	0.20%	1.72%
8	784 -> 784 -> 256	0.18%	1.68%

of the features extracted in the lower layers of the base classifiers are similar to one another, confirms the rationality of the mechanism. From Fig. 4(a), we see the weight increase of the difficult digits is rapid under boosting, which is the major reason that PB-DBN is achieving superior performance. And with more base classifiers involved in the committee judgement, the performance of PB-DBN improves (Table 2 and Fig. 4(b)).

## 4 Conclusion

In this paper, we described the Pseudo Boost Deep Belief Network(PB-DBN). Its unique strength is fast convergence to optimality while maintaining the same time complexity as DBN. In the model, top layers of DBNs are boosted based on different weighted data which helps the base classifiers(DBNs) to focus on more difficult data points, while the lower layers of the belief networks share weights for feature extraction. The method is computationally efficient and can be readily used for practical applications.

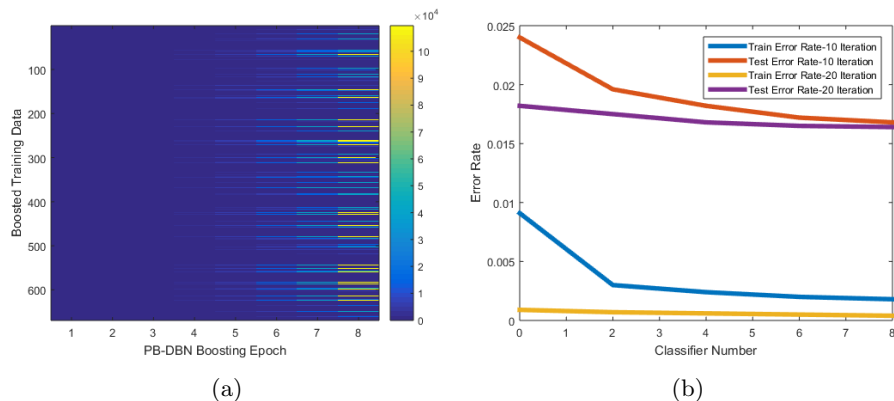
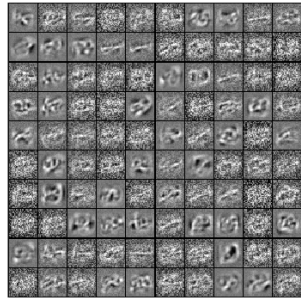


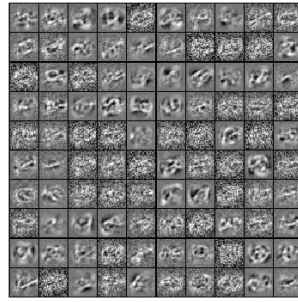
Fig. 4: (a) Evolution of boosted weights in PB-DBN. (b) Performance of PB-DBN with different number of base classifiers.

## References

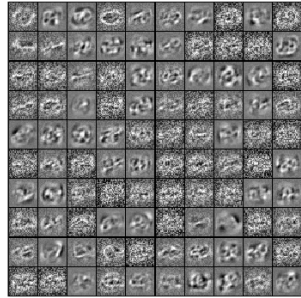
1. Andrew, G., Bilmes, J.: Sequential Deep Belief Networks. In: IEEE Int. Conf. Acoustics, Speech Signal Proc. (ICASSP). pp. 4265–4268 (2012)
2. Bengio, Y., Lamblin, P., Popovici, D., Larochelle, H., et al.: Greedy layer-wise training of deep networks. *Advances in neural information processing systems* 19, 153 (2007)
3. Bishop, C.M.: *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA (2006)
4. Bornschein, J., Bengio, Y.: Reweighted wake-sleep. CoRR abs/1406.2751 (2014)
5. Fischer, A., Igel, C.: An introduction to restricted Boltzmann machines. In: *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, pp. 14–36. Springer (2012)
6. Friedman, J., Hastie, T., Tibshirani, R.: Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *Ann. Statist.* 28(2), 337–407 (2000)
7. Hinton, G.E., Osindero, S., Teh, Y.W.: A fast learning algorithm for deep belief nets. *Neural Computation* 18(7), 1527–1554 (2006)
8. Lee, H., Ekanadham, C., Ng, A.Y.: Sparse deep belief net model for visual area v2. In: Platt, J.C., Koller, D., Singer, Y., Roweis, S.T. (eds.) *Advances in Neural Information Processing Systems 20*, pp. 873–880. Curran Assoc. (2008)
9. Lee, H., Grosse, R., Ranganath, R., Ng, A.Y.: Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In: *Proceedings of the 26th Annual International Conference on Machine Learning*. pp. 609–616. ICML '09, ACM, New York, NY, USA (2009)
10. Liu, P., Han, S., Meng, Z., Tong, Y.: Facial expression recognition via a boosted deep belief network. *IEEE Conf. Comput. Vision Pat. Recog. (CVPR)* (2014)
11. Liu, Y., Zhou, S., Chen, Q.: Discriminative deep belief networks for visual data classification. *Pattern Recognition* 44(10–11), 2287 – 2296 (2011)
12. Pham, T.V., Smeulders, A.W.M.: Quadratic boosting. *Pattern Recogn.* 41(1), 331–341 (Jan 2008)



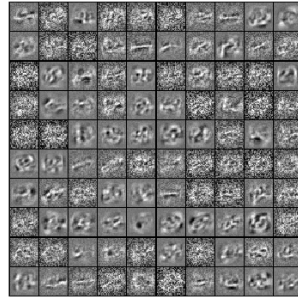
(a)



(b)



(c)



(d)

Fig. 5: Visualization of the feature extracted in the first hidden layer. (a) 2 layer DBN, 20 iterations training; (b) 3 layer DBN, 50 iterations training; (c) 2 layer PB-DBN, 20 iterations training; (d) 3 layer PB-DBN, 50 iterations training.

13. Salakhutdinov, R., Hinton, G.E.: Deep Boltzmann Machines. In: Proceedings of the International conference on artificial intelligence and statistics. vol. 5, pp. 448–455 (2009)
14. Sun, Y., Todorovic, S., Li, J.: Increasing the robustness of boosting algorithms within the linear-programming framework. *The Journal of VLSI Signal Processing Systems for Signal, Image, and Video Technology* 48(1), 5–20 (2007)
15. Wang, Y., Yang, J.A., Lu, J., Liu, H., Wang, L.W.: Hierarchical deep belief networks based point process model for keywords spotting in continuous speech. *International Journal of Communication Systems* 28(3), 483–496 (2015)
16. Zhiqiang, C., Chuan, L., Sánchez, R.V.: Multi-layer neural network with deep belief network for gearbox fault diagnosis. *Journal of Vibroengineering* 17(5), 2379 (2015)