Training Skinny Deep Neural Networks with Iterative Hard Thresholding Methods

Xiaojie Jin¹ Xiaotong Yuan⁴ Jiashi Feng² Shuicheng Yan^{3,2}

¹NUS Graduate School for Integrative Science and Engineering, NUS

²Department of ECE, NUS

³360 AI Institute

⁴School of Information and Control, Nanjing University of Information Science & Technology xiaojie.jin@nus.edu.sg, xtyuan1980@gmail.com, elefjia@nus.edu.sg, yanshuicheng@360.cn

Abstract

Deep neural networks have achieved remarkable success in a wide range of practical problems. However, due to the inherent large parameter space, deep models are notoriously prone to overfitting and difficult to be deployed in portable devices with limited memory. In this paper, we propose an iterative hard thresholding (IHT) approach to train Skinny Deep Neural Networks (SDNNs). An SDNN has much fewer parameters yet can achieve competitive or even better performance than its full CNN counterpart. More concretely, the IHT approach trains an SDNN through following two alternative phases: (I) perform hard thresholding to drop connections with small activations and fine-tune the other significant filters; (II) re-activate the frozen connections and train the entire network to improve its overall discriminative capability. We verify the superiority of SDNNs in terms of efficiency and classification performance on four benchmark object recognition datasets, including CIFAR-10, CIFAR-100, MNIST and ImageNet. Experimental results clearly demonstrate that IHT can be applied for training SDNN based on various CNN architectures such as NIN and AlexNet.

1 Introduction

Deep neural networks (DNNs) have achieved remarkable success in various applications. This has been driven by the rapid growth in the size of datasets, increasingly deeper network architectures and the development of various techniques in training deep models.

Despite their strong capability of learning rich and discriminative representations, DNNs usually suffer from following two problems caused by their inherent huge parameter space when applied in practice. First, most of state-of-the-art DNN architectures are prone to over-fitting even trained on large datasets [23, 28]. Secondly, as they usually consume large storage memory and computational resource, it is difficult to embed modern DNN models into devices with limited power and memory, e.g., mobile phones.

A lot of regularization techniques have thus been proposed to decrease the risk of over-fitting for DNNs, e.g., dropout [25]. However, those techniques are unable to reduce the storage cost. Recently, several works have been devoted to compressing and speeding-up DNNs by pruning internal layer connections. However, most of those methods achieve efficiency gain at the cost of performance deterioration.

In this work, we propose a novel approach for training skinny DNNs (SDNNs) with explicit size constraints to address the above problems simultaneously, *i.e.*, reducing the risk of overfitting to improve the generalization capability of deep models and meanwhile reducing the model size to decrease storage memory cost. Compared with other sparsity pursuit methods for training DNNs [9, 8], SDNN is superior as it is able to boost the performance even when the model

compression rate is high. Briefly, our proposed approach for training SDNNs contains two alternative phases:

- Phase I Hard thresholding over connections and sub-network fine-tuning. We apply hard thresholding over connections (weight parameters) at each layer to select the most prominent ones for the DNN model. The hard thresholding preserves the top k weight parameters with the largest magnitude and disables the others by zeroing their values. Then, we fine-tune the non-disabled parameters for compensating the performance loss caused by reducing the number of filters.
- **Phase II** Connection restoration and training the entire network. The frozen connections are re-activated and all the parameters are learned through training the entire network. The goal of this phase is to restore the truncated parameters and involve them in learning better representations.

Alternating the above two phases in training DNNs is able to produce SDNNs which have a stronger generalization capability with fewer parameters compared with the counterpart trained in the conventional way. We term such a method compositing of the above two phases as iteratively hard thresholding (ITH). Note that these two alternative phases are only needed in training SDNNs, while in the testing stage SDNN only takes one feedforward pass for inputs to make prediction.

To verify the effectiveness of SDNNs and the ITH training approach, we conduct extensive experiments on four public datasets with various scales, *i.e.*, MNIST [17], CIFAR10 [14], CIFAR100 [14] and ImageNet [4] for two DNN architectures with different complexities including Network in Network [21], and AlexNet [15]. The experimental results clearly demonstrate that SDNN with ITH does not only improve the generalization capability of deep models and provide state-of-the-art performance, but also reduce the size of parameters at the same time. Therefore, ITH is a quite appealing approach for training SDNNs in real-world scenarios concerning limited computational and storage cost.

2 Preliminaries: Gradient Hard Thresholding Revisit

The gradient hard thresholding (GHT) algorithm was proposed by [30] to solve the following sparsity-constrained convex optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}), \text{ s.t. } \|\mathbf{x}\|_0 \le k, \tag{1}$$

where $f: \mathbb{R}^d \to \mathbb{R}$ is a smooth convex function and $\|\mathbf{x}\|_0$ counts nonzero elements in \mathbf{x} . The GHT algorithm solves the problem in Eqn. (1) by alternatively performing gradient descent and hard thresholding over the gradient. More concretely, let $\mathcal{O}_k(\mathbf{x})$ denote the hard thresholding operator which selects the top k entries of \mathbf{x} with largest magnitudes and set the rest entries to 0. Let $\mathbf{x}^{(t)}$ denote the updated variable and $F^{(t)}$ denote the support set of $\mathcal{O}_k(\mathbf{x}^{(t)})$ at the t-th iteration. We further define supp(\mathbf{x}) as the support set of \mathbf{x} and supp(\mathbf{x} , k) as the index set of the top k entries of \mathbf{x} . At the t-th iteration, following three steps (denoted as S1, S2 and S3, respectively) are involved in GHT:

S1: perform gradient descent at $\mathbf{x}^{(t-1)}$ with a step-size η : $\tilde{\mathbf{x}}^{(t)} = \mathbf{x}^{(t-1)} - \eta \nabla f(\mathbf{x}^{(t-1)})$, where $\nabla f(\mathbf{x}^{(t-1)})$ is the gradient of $f(\cdot)$ evaluated at $\mathbf{x}^{(t-1)}$;

S2: apply hard thresholding $\mathcal{O}_k(\cdot)$ on $\tilde{\mathbf{x}}^{(t)}$ as: $\tilde{\mathbf{x}}^{(t)} = \mathcal{O}_k(\tilde{\mathbf{x}}^{(t)})$. Therefore, $F^{(t)} = \operatorname{supp}(\tilde{\mathbf{x}}, k)$; S3: optimize $\mathbf{x}^{(t)}$ by minimizing the objective function over support set $F^{(t)}$, i.e., $\mathbf{x}^{(t)} = \arg\min\{f(\tilde{\mathbf{x}}), \sup_{t \in \mathcal{X}} \mathbf{x}^{(t)}\}$.

It is proved that under mild conditions, the GHT algorithm converges geometrically to the point with bounded deviation from global optimum, with a high probability [30].

```
Algorithm 1: Iterative Hard Thresholding for Training SDNNs
```

```
input : Training sample (\mathbf{x}, y^*), randomly initialized weights \mathbf{W} = (\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}) for an L-layer SDNN. Parameters s_1 and s_2. Loss function \mathcal{L}.

(Step 1) Train the SDNN model for s_1 epochs, i.e., \mathbf{W} = \arg\min_{\mathbf{W}} \mathcal{L}(\mathbf{x}, y^*, \mathbf{W}). while maximum number of epochs is not reached do

Phase I:

(Step 2) Apply hard thresholding operation \mathcal{O}_{k_\ell}(\mathbf{x}) to \mathbf{W}^{(\ell)}: \mathbf{W}^{(\ell)} = \mathcal{O}_{k_\ell}(\mathbf{W}^{(\ell)}).

(Step 3) F^{(\ell)} = \sup_{\mathbf{W}}(\mathbf{W}^{(\ell)}, k_\ell), \ \ell = 1, \dots, L.

(Step 4) Train the SDNN for s_2 epochs:

\mathbf{W} = \arg\min_{\mathbf{W}} \{\mathcal{L}(\mathbf{x}, y^*, \mathbf{W}), \sup_{\mathbf{W}^{(\ell)}} \subseteq F^{(\ell)} \}.

Phase II:

(Step 5) Restore the connections truncated at Step 2.

(Step 6) Train the deep model for s_1 epochs: \mathbf{W} = \arg\min_{\mathbf{W}} \{\mathcal{L}(\mathbf{x}, y^*, \mathbf{W})\}.
```

3 Skinny Deep Neural Networks

3.1 Deep Neural Networks with Cardinality Constraint

For expression conciseness, we only consider the case with a single training sample. The formulation for training with multiple samples can be derived similarly as samples are independent and the loss function is decomposable over samples. We denote a training sample as (\mathbf{x}, y^*) where $\mathbf{x} \in \mathbb{R}^d$ denotes the raw input data and $y^* \in \{1, \ldots, C\}$ is its ground truth category label. Here C is the total number of categories. We consider a DNN model consisting of L layers, each of which outputs a feature map, denoted as $\mathbf{X}^{(\ell)}$ for layer ℓ . Here $\mathbf{X}^{(0)}$ and $\mathbf{X}^{(L)}$ represent the input and final output of the network, respectively. Let $\mathbf{W}^{(\ell)}$ denote parameters of the filters (or weights) of the ℓ -th layer to be learned, and $\mathbf{W} = (\mathbf{W}^{(1)}, \ldots, \mathbf{W}^{(L)})$ is the collection of all learnable parameters in a DNN. We use $\mathbf{W}^* = (\mathbf{W}^{(1)^*}, \ldots, \mathbf{W}^{(L)^*})$ to denote output the parameters after training. Using the above notations, the output of each layer in an L-layer DNN can be written as

$$\mathbf{X}^{(\ell)} = g^{(\ell)}(\mathbf{W}^{(\ell)} * \mathbf{X}^{(\ell-1)}), \quad \ell = 1, \dots, L \quad \text{and} \quad \mathbf{X}^{(0)} \triangleq \mathbf{x},$$

where $g^{(\ell)}(\cdot)$ is a composite of multiple specific functions including activation function, dropout, pooling, batch normalization and softmax. For succinct notations, the bias term is omitted. The loss function of a deep model that we consider here is

$$\mathcal{L}(\mathbf{x}, y^*, \mathbf{W}) = -\log(h_{v^*}(\mathbf{x}, \mathbf{W})) + \lambda \|\mathbf{W}\|_F, \tag{2}$$

where $h_{y^*}(\mathbf{x}, \mathbf{W})$ denotes the probability score predicted for \mathbf{x} on the ground truth category of y^* , and λ is the weight decay factor. In traditional methods, a deep model is optimized through minimizing the loss function without any constraint over \mathbf{W} :

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{x}, y^*, \mathbf{W}). \tag{3}$$

In this work, we propose to impose explicit cardinality constraints to layer-wise parameters during training in order to reduce the parameter size. Thus the optimization problem for training an SDNN is formulated as

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{x}, y^*, \mathbf{W}), \text{ s.t. } \|\mathbf{W}^{(\ell)^*}\|_0 \le k_{\ell}, \quad \ell = 1, \dots, L,$$

$$\tag{4}$$

where k_{ℓ} is the cardinality constraint for the parameters of the ℓ -th layer. Note that k_{ℓ} can be either the same or different across different layers. Correspondingly, $r^{(\ell)}$ is denoted as the sparsity ratio of the ℓ -th layer and defined as $r^{(\ell)} \triangleq \|\mathbf{W}^{(\ell)}\|_0 / |\mathbf{W}^{(\ell)}|$ where $|\mathbf{W}^{(\ell)}|$ denotes the number of parameters at the ℓ -th layer.

3.2 From GHT to SDNN

Inspired by the success of GHT on solving sparsity-constrained convex problems [30], we apply it to train DNNs to alleviate the over-parameterization issue. However, straightforwardly applying GHT to train SDNNs cannot provide desirable results since the loss function is non-convex and highly complex. Therefore applying hard thresholding operation to the parameters of deep model at each iteration in the same way as GHT will cause the model to diverge and be unable to learn meaningful parameters.

In this paper, we propose a novel approach for training SDNNs with cardinality constraints. A summary of details on training SDNNs is presented in Algorithm 1. Compared with GHT, there are two main differences in the approach used for training SDNNs:

Multi-Step Update. GHT performs hard thresholding operation at each iteration. However in training SDNNs, we update the parameters for s_1 iterations before performing hard thresholding. Such a strategy yields following two advantages. First, the training is accelerated by largely reducing the time of performing hard thresholding. Secondly, by updating all parameters (including those of connections truncated by the last hard thresholding operation) sufficiently, the SDNNs are likely to learn more discriminative representation with more parameters. Otherwise with a single-step update, it is highly possible that connections which are truncated by the first hard thresholding operation would be always truncated in all subsequent hard thresholding operations due to the little change in magnitudes.

Relaxation on Sparsity Constraints. As can be seen from Eqn. (1), GHT applies the ℓ_0 sparsity constraints to the parameters at each iteration. In contrast, since we only care about the final model after training is finished, we do not need to optimize the SDNNs with sparsity constraints throughout the training stage. Actually, in SDNNs, we use such relaxation to restore the connections truncated by hard thresholding operation, aiming to learn better representations. More details are given in Section 3.3. As confirmed by experimental results, it is critical for SDNNs to learn more discriminative features than other network pruning works [9, 8, 2].

3.3 Training SDNNs

The training of SDNN mainly consists of two iteratively alternating phases, each of which contains an operation for the parameters of deep models followed by an optimization process which is mutually different. In this section, we explain the details of Algorithm 1.

Network Initialization: This corresponds to Step 1. At the beginning of the training stage, all the parameters of the SDNN are trained for s_1 epochs without considering cardinality constraints. Therefore the optimization formulation in this step is exactly the same as the one in Eqn. (3). The aim of this step is to provide a good initialization for the following steps, which prevents the SDNN from diverging or getting stuck in a bad local minimum.

Phase I: This phase corresponds to Step 2, Step 3 and Step 4 in Algorithm 1, which perform hard thresholding and sub-network fine-tuning, respectively. At Step 3, we keep the top k_{ℓ} parameters with largest magnitudes at the ℓ -th layer and set the rest parameters to 0. At Step 4, we fine-tune the parameters reserved at Step 3. This step aims to compensate the performance drop caused by erasing a part of the parameters. As demonstrated in our experiments, reinitializing the reserved parameters leads to poor performance. This phenomenon could be explained similar as the observations made in [29]: since DNNs contain fragile coadapted features, the gradient descent algorithm is able to find a good solution when the network is trained from the scratch, but it may fail after re-initializing some layers and retraining them.

Moreover, starting to fine-tune the network with retained weights requires less computation because it is not necessary to train the entire network.

Phase II: This phase corresponds to Step 5 and Step 6, which conduct a weight restoration operation and trains the entire network, respectively. Step 5 removes the cardinality constraints over parameters which are set to be 0 in Step 3 so that all parameters are updated in Step 6 freely. This step is critical in SDNN for the following reason: by updating all parameters including those set to be 0 in Step 3, the SDNN is able to restore some connections that are beneficial for learning feature representation with strong discriminative power. As can be seen in Figure 2, at the initial training stage, a large proportion of connections which are truncated in the last round of hard thresholding operation become significant (with parameters in large magnitudes) at this phase. Thus SDNN has a strong capability to search among a large parameter space for seeking a better local optimum. Besides, we observe that the ratio decrease as the training progresses due to the deep model converges to a good optimum.

4 Related Works

Early approaches for deep model compression including optimal brain damage [18] and optimal brain surgeon [10] that prune the connections in networks based on the second order information. However, those methods are not feasible for deep networks due to high computational complexity. Recent works aiming at network pruning include [9, 3, 2, 16, 13], which prune connections in a progressively greedy way [9] or using sparsity related regularizer [3, 16]. Although those works can reduce model size significantly, they suffer from the dramatic performance loss. In contrast, SDNN does not only offer significant compression ratios but also improves the performance simultaneously.

Our work is also in line with model compression. For example, [2] proposes to quantize the deep model by minimizing L2 error and [5] seeks an low-rank approximation of the model. Recently, [8] combined pruning, quantization and Huffman coding techniques and provided rather high compression ratios. However, those methods also introduce performance drop. There are also works trying to compress a model by using binning methods [6], but they can only be applied over fully connected layers. In contrast, our method can be applied for compressing both convolution layers and fully connected layers.

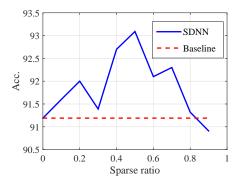
5 Experiments and Analysis

5.1 Experiment Settings and Implementation Details

We justify our method on four scale-various object classification benchmarks, *i.e.* three small-scale ones including CIFAR10 [14], CIFAR100 [14] and MNIST [17] and the large-scale ImageNet [4]. Two evaluation metrics, the classification performance and the number of parameters in a model, are used in comparison with other methods.

Deep Models In our experiments, we train and test two deep models which are with different complexities, *i.e.* Network in Network (NIN) [21], and AlexNet [15]. Briefly, NIN has only convolutional weight layers by replacing the single linear convolution layers in the conventional CNNs by multilayer perceptrons and using the global average pooling layer to generate feature maps for each category. Compared with NIN, AlexNet is wider and deeper containing 60M parameters with five convolution layers and three fully connected layers.

Implementation All of our experiments are conducted on a NVIDIA TITAN GPU using Caffe. The numbers of training epochs for Step 1 and Step 4 in Algorithm 1 are set to $s_1 = 15$ and $s_2 = 150$ for small-scale datasets and $s_1 = 15$ and $s_2 = 40$ for the ImageNet dataset to reduce training time. The hyperparameters of NIN and AlexNet including learning rate, momentum



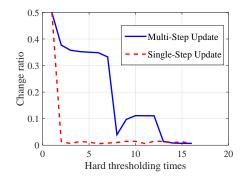


Figure 1: The accuracy of SDNN on CI-FAR10 with different sparse ratios. The baseline model is equal to a SDNN with sparse ratio r=0. From r=0.1 to r=0.8, SDNN surpasses the baseline model while reducing the model size. Particular, when r=0.5, SDNN significantly improves the baseline by 1.9%. When r=0.9, SDNN is outperformed by baseline model by only 0.28%.

Figure 2: The change ratio of connection truncated by neighboring hard thresholding operations. The multi-step update adopted by SDNN has large change ratio in initial training phase, allowing the deep model to learn features with strong discriminative capability, while the single-step update would result $\sim 0.01\%$ change ratios during all training phase, thus undermines the representation capability of deep models.

and weight decay follow [21] and [15], respectively. All the results of our methods in the paper are based on the models reaching convergence.

Data augmentation is used in many models for object classification to prevent overfitting. The horizontal flipping is used for CIFAR10 and CIFAR100. For ImageNet, we use random crop and horizontal flipping as in [15]. No data augmentation method is applied for MNIST.

Memory Usage To efficiently utilize the sparsity property of SDNN to reduce the memory storage size, we refer to the "Bitmask" storage format proposed in [3], which stores the nonzero parameters as well as a mask whose number of bits is equal to the number of total parameters. The bit value will be set to 1 if the corresponding parameter is nonzero, otherwise to 0. As indicated by [3], such memory usage methods can be directly used in deep models at runtime and reduce the storage size efficiently at the same time.

5.2 Model Analysis

In the following, we use NIN trained and tested on CIFAR10 to investigate the effects of different k_{ℓ} , $\ell = 1, ..., L$ (recall that k_{ℓ} denotes the number of nonzero parameters in the ℓ -th layer) on the classification performance of SDNN as well as justify the hard thresholding strategy by comparing with the stochastic thresholding strategy. As a baseline, the NIN trained without sparsity constraint achieved 91.9% [21] with data augmentation on this dataset.

Global Sparsity Distribution One problem in practical use is how to set the sparsity constraint $k_{(\ell)}$ for different layers. To ease the laborious work of tuning such parameters, especially when the number of layers in a deep model is large, e.g. ResNet [11], we explore a simple way by setting all layers in a deep model with the same sparsity constraints. As we have observed in the experiment, the training process of the network may diverge if directly applying hard thresholding to the layer-wise parameters when $r^{(\ell)}$ is large, e.g., 0.8. We propose a progressive hard thresholding strategy to address this issue, i.e., $r^{(\ell)}(t) = r^{(\ell)}(0) + t(r^{(\ell)}(T) - r^{(\ell)}(0))/T$ where T is the overall training epochs, $r^{(\ell)}(t)$ denotes the sparsity ratio to layer-wise parameters when applying hard thresholding operation at t-th epochs. $r^{(\ell)}(0)$ and $r^{(\ell)}(T)$ are the sparsity ratio when the first hard thresholding operation is applied and the target layer-wise sparsity ratio, respectively. The experimental results prove that such a simple strategy is able to boost

the performance of the baseline model and simultaneously reduce the model size. Note although Collins et~al.~[3] proposed an iterative method to search the layer-wise sparsity hyperparameters in a greedy way, their method is time consuming thus not feasible for large deep models. In the experiments, we set the same sparsity ratio $r = \{0, 0.1, 0.2, \ldots, 0.9\}$ for all layers. Note when r = 0, our model degenerates to the baseline model. Figure 1 illustrates the curve of classification performance versus the sparsity ratio. As obviously observed in Figure 1, SDNN outperforms the baseline model with the sparsity ratio ranging from 0.1 to 0.8, which demonstrates the strong capability of SDNN in enhancing the generalization capability of deep models. Particularly, when r = 0.5, SDNN significantly increases the classification accuracy by 1.90%. Even when the largest ratio is increased to r = 0.9, i.e. reducing the model size by $10 \times$, SDNN only bears a slight performance loss (0.28%) compared to the baseline model. Note that in [2], the accuracy of CIFAR10 using their method started to drop from r = 0.6, which was much earlier than ours. Above experimental justify the effectiveness of applying IHT to train SDNN.

Hard Thresholding vs Random Thresholding We conduct experiments to justify the hard thresholding by replacing it with random thresholding in our method while keeping the other configurations unchanged. Specifically, we set r=0.5 in our experiments and replace Step 4 in Algorithm 1 with random thresholding in which half of parameters in each layer are randomly set to zero. However, we observe that deep models diverge quickly in Step 5 after random thresholding. The reason is that the random thresholding places deep models in bad conditions by truncating those important parameters.

5.3 Results

CIFAR10 The CIFAR-10 dataset consists of 60,000 color images of 32×32 pixels in 10 classes. The total dataset is split into 50,000 training images and 10,000 testing images. Table 1 compares the performance and # parameters of SDNN and other state-of-the-art methods either when data augmentation is used or not. SDNN with different sparsity ratios are denoted with SDNN-#× where # is the reciprocal of the sparsity ratio. It is observed that when data augmentation is not applied, SDNN with sparsity ratio r = 0.5 achieves the best result among all methods, reducing the error rate (ER) of the baseline model NIN by 1.71%. Compared with RCNN-96 and RCNN-128 which have model sizes of 0.67M and 1.19M, respectively, SDNN-2× with a much smaller model size (0.49M) outperforms them by 0.61% and 0.28%, respectively, demonstrating that SDNN is able to significantly improve the generalization capability of deep models. To further test SDNN's performance in deep models with a larger size, we evenly increase the parameter of each layer in the original NIN by two times, resulting in a model which is four times as large as the original NIN. We denote the enlarged NIN as NIN2 and test SDNN with sparsity ratio r = 0.5 (denoted as SDNN₂-2×) on it. Compared with NIN₂, SDNN₂-2× is able to reduce the ER by 0.77% and 1.72% when data augmentation is applied or not, respectively, again verifies SDNN's capability to reduce overfitting. Note compared with ResNet which gets the lowest ER against all other methods, our method is only with a 0.02% higher ER, but our model is much faster in testing since ResNet with 1.7M parameters has 110 layers while ours only has 9 layers.

We also test SDNN with larger sparsity ratios when data augmentation is used. When r = 0.9, SDNN-10x is slightly worse than the baseline NIN model (0.28% higher on ER) with a significantly reduced model size (0.1M). We further train the SDNN-20x which has a sparsity ratio r = 0.95 with model size 0.05M to achieve the ER of 10.53%, which is 1.64% higher than the baseline NIN. Considering NIN is a network consisting of all convolution layers, this is still a satisfiable result.

Table 1: Comparison with existing models on CIFAR10. ER represents "error rates" and DA represents "data augmentation". The subscript for NIN and SDNN represents the ratio of the parameters' number in each layer w.r.t. the model without subscript. E.g., NIN₂ has as two times as many layer-wise parameters compared to baseline NIN. The number in " $\#\times$ " for SDNN represents the model compression ratio of that model. The number after RCNN represents the number of convolution feature maps in each layer. Best results are shown in bold.

Model	No. of Param.(M)	ER w/ DA(%)	ER w/o DA(%)
NIN [21]	0.97	10.41	8.81
DSN [19]	0.97	9.69	7.97
RCNN-96 [20]	0.67	9.31	7.37
RCNN-128 [20]	1.19	8.98	7.24
RCNN-160 [20]	1.86	8.69	7.09
FitNet [22]	~ 2.5	-	8.39
Highway [26]	2.3	-	7.54
ResNet [11]	0.27	-	8.75
ResNet [11]	1.7	-	6.43
$NIN_{1/2}$	0.49	10.50	9.20
NIN_2	3.92	9.2	8.17
SDNN-2×	0.49	8.70	6.91
SDNN-10 \times	0.1	-	9.09
SDNN-20 \times	0.05	-	10.53
$SDNN_2$ -2×	1.79	8.43	$\boldsymbol{6.45}$

Table 2: Comparison with existing models on CIFAR100. Note since NIN [21] did not report the results with DA, we therefore refer to the reimplementation of NIN in [1] when DA is used.

Model	No. of Param.(M)	ER w/ DA(%)	ER w/o DA(%)
Maxout [7]	>5	38.57	-
Prob. maxout [24]	>5	38.14	-
NIN [21]	>5	35.68	-
NIN [1]	0.97	35.96	32.70
DSN [19]	0.97	34.57	-
dasNet [27]	-	34.50	-
RCNN-96 [20]	0.67	34.18	-
RCNN-128 [20]	1.19	32.59	-
RCNN-160 [20]	1.86	31.75	-
APL [1]	0.97	34.40	30.83
SReLU [12]	0.97	31.23	29.91
NIN_2	3.81	31.29	29.12
SDNN-2×	0.49	30.50	29.51
$SDNN_2-2\times$	1.79	29.13	27.14

CIFAR100 CIFAR100 has 50,000/10,000 training/testing color images with resolution of 32×32 pixels. Since the number of training images of each class in CIFAR100 is only one tenth of that in CIFAR10, deep models trained on this dataset are prone to overfitting. Table 2 summaries the state-of-the-art methods on this dataset. It is obviously observed that either when the data augmentation is used or not, SDNN-2× achieves the lowest ER among all methods with the smallest model size. Specifically, compared with NIN model, SDNN-2× is able to reduce the ER by 5.18% and 3.19% either when data augmentation is used or not, respectively, which again demonstrates the advantages of SDNN. Moreover, compared with RCNN-96 which has a model size of 0.67M and achieves 34.18% ER without using data augmentation, our method outperforms it significantly by reducing the ER by 3.68% while simultaneously attaining a smaller model size (0.49M versus 0.67M). Like on CIFAR10, we also test SDNN₂-2× by adopting half of the model size of NIN₂. As can be seen from Table 2, our method is able to further enhance the generalization capability of NIN₂ by reducing the ER significantly by 2.16% /1.98% when trained with / without data augmentation, respectively.

Table 3: Error rates on MNIST without data augmentation.

-	.0		
	Model	# Param.(M)	ER(%)
	Stocha. Pooling [31]	-	0.47
	Maxout [7]	0.42	0.45
	NIN [21]	0.35	0.47
	DSN [19]	0.35	0.39
	SReLU [12]	0.35	0.35
	RCNN-96 [20]	0.67	0.31
	SDNN-2×	0.18	0.19

Table 4: Comparison with existing models on ImageNet.

Model	# Param.(M)	Top-5 $ER(\%)$
Caffe Version ¹	60	19.56
Han & Pool [9]	6.7	19.67
Mem. Bound [3]	15	19.6
SDNN-2×	30	17.90
SDNN-4 \times	15	18.75

MNIST MNIST is one of the standard datasets in the machine learning community. It consists of 70,000 handwritten digits of 0 to 9 with 28×28 resolution of pixels in gray scale format, which are split into 60,000/10,000 training/testing set. Table 3 shows the comparison results, from which we can see SDNN-2× is able to achieve the best performance using the least parameters. Compared with the baseline model NIN, SDNN-2× greatly reduces the error rates from 0.47% to 0.19% on this heavily benchmarked dataset with only half of the NIN model size. To our best knowledge, this is the best performance ever achieved by methods without using ensembles and other preprocessing methods, which could further boost the performance of ours.

ImageNet To test the scalability of SDNN to large-scale datasets and model with bigger sizes, we conduct experiments on a much more challenging image classification task on 1000-class ImageNet dataset. As a renowned dataset in the vision research community, ImageNet contains about 1.2M training images, 50,000 validation images and 10,0000 testing images. To compare with other network pruning methods [9, 3], we also use AlexNet as our baseline model according to the publicly available configurations in Caffe¹. Table 4 lists the performance of SDNN and other methods on ImageNet using AlexNet. Compared with the baseline model, our model with a two and four times compression ratio can further reduce the top-5 error rates by 1.66% and 0.81%, respectively, demonstrating the efficacy of SDNN on large-scale datasets. Compared with [9, 3], both report higher error rates when pruning the deep model. SDNN is superior by boosting the performance of deep models while largely reducing the model size at the same time. For example, compared with [3], SDNN-4× reduces its error rates significantly by 0.85% while obtaining models with the same size.

6 Conclusion

In this paper, we proposed an iterative hard thresholding method (IHT) to improve the performance of the deep neural network and reduce the size of parameters simultaneously. The training of SDNN using IHT consists of two alternative phases, *i.e.* firstly performing hard thresholding to set connections with small magnitudes to zero and fine-tune the significant filters, and secondly, re-activating the freezing. Experiments conducted on four scale-various datasets, *i.e.* CIFAR10, CIFAR100, MNIST and ImageNet using deep networks with different complexities, *i.e.* NIN and AlexNet demonstrated that our method is able to significantly boost the discriminative capability of deep models while largely reducing their sizes simultaneously.

References

[1] Forest Agostinelli, Matthew Hoffman, Peter J. Sadowski, and Pierre Baldi. Learning activation functions to improve deep neural networks. *CoRR*, abs/1412.6830, 2014.

¹https://github.com/BVLC/caffe/tree/master/models/bvlc-alexnet

- [2] Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural networks. arXiv preprint arXiv:1512.08571, 2015.
- [3] Maxwell D Collins and Pushmeet Kohli. Memory bounded deep convolutional networks. arXiv preprint arXiv:1412.1442, 2014.
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [5] Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in Neural Information Processing Systems*, pages 1269–1277, 2014.
- [6] Yunchao Gong, Liu Liu, Ming Yang, and Lubomir Bourdev. Compressing deep convolutional networks using vector quantization. arXiv preprint arXiv:1412.6115, 2014.
- [7] Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. *arXiv preprint arXiv:1302.4389*, 2013.
- [8] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. CoRR, abs/1510.00149, 2015.
- [9] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In Advances in Neural Information Processing Systems, pages 1135–1143, 2015.
- [10] Babak Hassibi and David G Stork. Second order derivatives for network pruning: Optimal brain surgeon. Morgan Kaufmann, 1993.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.
- [12] Xiaojie Jin, Chunyan Xu, Jiashi Feng, Yunchao Wei, Junjun Xiong, and Shuicheng Yan. Deep learning with s-shaped rectified linear activation units. *CoRR*, abs/1512.07030, 2015.
- [13] Yong-Deok Kim, Eunhyeok Park, Sungjoo Yoo, Taelim Choi, Lu Yang, and Dongjun Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. arXiv preprint arXiv:1511.06530, 2015.
- [14] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images, 2009.
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.
- [16] Vadim Lebedev and Victor Lempitsky. Fast convnets using group-wise brain damage. arXiv preprint arXiv:1506.02515, 2015.
- [17] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [18] Yann LeCun, John S Denker, Sara A Solla, Richard E Howard, and Lawrence D Jackel. Optimal brain damage. In *NIPs*, volume 89, 1989.
- [19] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. arXiv preprint arXiv:1409.5185, 2014.

- [20] Ming Liang, Xiaolin Hu, and Bo Zhang. Convolutional neural networks with intra-layer recurrent connections for scene labeling. In *NIPS*, 2015.
- [21] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. arXiv preprint arXiv:1312.4400, 2013.
- [22] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550, 2014.
- [23] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.
- [24] Jost Tobias Springenberg and Martin Riedmiller. Improving deep neural networks with probabilistic maxout units. arXiv preprint arXiv:1312.6116, 2013.
- [25] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [26] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. *CoRR*, abs/1505.00387, 2015.
- [27] Marijn F. Stollenga, Jonathan Masci, Faustino J. Gomez, and Jürgen Schmidhuber. Deep networks with internal selective attention through feedback connections. CoRR, abs/1407.3068, 2014.
- [28] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. arXiv preprint arXiv:1409.4842, 2014.
- [29] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, pages 3320–3328, 2014.
- [30] Xiao-Tong Yuan, Ping Li, and Tong Zhang. Gradient hard thresholding pursuit for sparsity-constrained optimization. arXiv preprint arXiv:1311.5750, 2013.
- [31] Matthew D. Zeiler and Rob Fergus. Stochastic pooling for regularization of deep convolutional neural networks. CoRR, abs/1301.3557, 2013.