

# Recent Advances in Large Margin Learning

Yiwen Guo and Changshui Zhang, *Fellow, IEEE*

**Abstract**—This paper serves as a survey of recent advances in large margin training and its theoretical foundations, mostly for (nonlinear) deep neural networks (DNNs) that are probably the most prominent machine learning models for large-scale data in the community over the past decade. We generalize the formulation of classification margins from classical research to latest DNNs, summarize theoretical connections between the margin, network generalization, and robustness, and introduce recent efforts in enlarging the margins for DNNs comprehensively. Since the viewpoint of different methods is discrepant, we categorize them into groups for ease of comparison and discussion in the paper. Hopefully, our discussions and overview inspire new research work in the community that aim to improve the performance of DNNs, and we also point to directions where the large margin principle can be verified to provide theoretical evidence why certain regularizations for DNNs function well in practice. We managed to shorten the paper such that the crucial spirit of large margin learning and related methods are better emphasized.

**Index Terms**—Large margin classifier, adversarial perturbation, generalization ability, deep neural networks

## 1 INTRODUCTION

THE concept of large margin learning arises along with the development of support vector machine (SVM) [1], [2], which aims to fix the empirical risk and minimize the confidence interval, in contrast to many models that target mostly at minimizing the empirical risk [2], [3]. Benefit from solid theoretical basis from statistical learning, large margin classifiers show promise in both generalization ability and robustness. Since the late 1990s, they have been intensively studied and widely adopted [4], [5], [6], [7], [8].

Yet, every learning machine has its day. Recent years have witnessed a revive of neural networks, partially owing to the advances of computational units that are capable of processing large-scale datasets. By learning representations from data, they, especially deep neural networks (DNNs), have advanced the state-of-the-arts of many tasks for machine intelligence [9], [10]. One might be curious about the relationship between DNNs and conventional large margin classifiers (e.g., SVM), and in view of this, we would like to answer three questions in this survey: **1)** Is the large margin principle essential or at least beneficial to the classification performance of DNNs? **2)** if yes, is it (implicitly) supported with a normal training mechanism used in practice? **3)** If not implicitly supported, how to gain large margins for classification models that are nonlinear and structural complex like DNNs? This paper presents an overview of existing work on these points. To the best of our knowledge, there is no such survey in the literature, and our work will bridge the information gap and inspire new research hopefully.

### 1.1 A Formal Definition of Classification Margin

To get started, let us first introduce a formal definition of the **classification margin**, which applies to a variety of different classifiers, including both linear and nonlinear ones.

- Y. Guo is with ByteDance AI Lab, Beijing 100000, China. E-mail: guoyiwen.ai@bytedance.com.
- C. Zhang is with the Institute for Artificial Intelligence, Tsinghua University (THUAI), Beijing National Research Center for Information Science and Technology (BNRist), and the Department of Automation, Tsinghua University, Beijing 100084, China. E-mail: zcs@mail.tsinghua.edu.cn

Manuscript received Oct. 20, 2019, revised Mar. 20, 2021.

It can be slightly different to formulate the margin of binary and multi-class classification. We first consider the binary scenario, in which a label  $y$  from  $\{+1, -1\}$  is assigned by a classifier to its input  $\mathbf{x} \in \mathbb{R}^n$ . Given  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  which maps the  $n$ -dimensional input  $\mathbf{x}$  into a one-dimensional decision space, an *instance-specific margin* is defined as

$$m_{\mathbf{x},f} := \min \|\mathbf{z}\|_p \text{ s.t. } f(\mathbf{x})f(\mathbf{x} + \mathbf{z}) < 0 \quad (1)$$

for any  $\mathbf{x} \in \mathbb{R}^n$ .  $p \geq 1$  indicates the concerned  $l_p$ -norm, and, in the Euclidean space, we have  $p = 2$ . For a classifier whose prediction is given by a linear function of its input  $\mathbf{x}$ , we can rewrite the function output  $f(\mathbf{x})$  as  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ , and the instance-specific margin has a closed-form solution in this case:  $m_{\mathbf{x},f} = |f(\mathbf{x})|/\|\mathbf{w}\|$ . It can be of great interest to study a margin  $m_f := \min m_{\mathbf{x},f}$  over the whole training set. SVM in the linear case enlarges such a classification margin by minimizing  $\|\mathbf{w}\|$  and constraining the value of  $f(\mathbf{x})$  (to be more specific, by letting  $yf(\mathbf{x}) > 1$  for all  $\mathbf{x} \in \mathbb{R}^n$ ). Taking advantage of some kernel trick [11] and utilizing nonlinear functions in reproducing kernel Hilbert spaces, it is natural to extend the expressions and analyses of margins to SVMs with higher expressivity [3].

With the expression of classification margin provided for SVMs, we obtain certification on the superiority of their generalization ability and robustness [12], [13], [14]. For DNNs, it is very challenging, if not impossible, to derive analytic expressions for their classification margins, on account of the hierarchical model structure and complex nonlinearity of  $f$ . Therefore, the study of classification margin, generalization ability, and robustness of DNNs also attract great attention recently. In the following sections, we will try to answer the three questions and given an overview of recent advances in margin related studies for DNNs. Before providing more details, it is worth stressing that the “classification margin” defined in the beginning of this section and concerned in linear SVMs is in the *input space* of learning models. We will also mention some margins in the output or representation spaces of DNNs in this paper.

## 2 THE RELATIONSHIPS BETWEEN MARGIN, GENERALIZATION, AND ROBUSTNESS

In this section, we attempt to answer the first two questions raised in Section 1. First, it is the question about benefits of large margin learning to DNN-based classifications. Related studies from theoretical perspectives have been performed for decades on the basis of some shallow models like SVMs. In this section, we focus on research work based on DNNs.

### 2.1 Large Margin Is Beneficial to DNNs as Well

Before delving deep into these studies, we first introduce the concept of **robustness**, (a.k.a., algorithmic robustness [15] and robustness in patterns [8]), which is an essential property of classifiers and closely related to the classification margin. As is known from the definition, for any reasonable input (e.g., any natural image that can be fed into a scene classification system), the classifier will hold its prediction with any pixel-wise perturbation smaller than the margin (i.e.,  $\|z\|_p < m_f$ ). We consider  $f$  as a robust model if  $m_f$  is sufficiently large so that the perturbation lead to samples perceptually from the other classes. Over the last few years, the robustness of DNN models has draw more attention along with developments of adversarial attacks [16], [17], [18], [19], [20]. It has been demonstrated that one can easily manipulate the prediction of a state-of-the-art DNN model by adding subtle perturbations to its input. By definition, given the function  $f$  and an input  $x$ , any perturbation to  $x$  within the hyper-sphere  $\{z \mid \|z\| < m_{x,f}\}$  would not alter the model prediction. That said, the concept of adversarial robustness that describes the ability of a model to resist adversarial attacks, is intrinsically related to the classification margin.

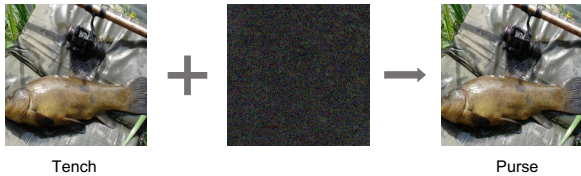


Fig. 1. An example of the adversarial example which is misclassified as purse by a ResNet-50 trained on ImageNet. We enlarge the perturbation by  $10\times$  for better illustration in the middle picture.

What relates to robustness and being important as well is the **generalization ability**. Suppose that a DNN model is to be learned to minimize the risk  $\tilde{l}(\Theta) := \mathbb{E}_{\mathbf{x},y}(l(f(\mathbf{x}; \Theta), y))$ , in which  $\Theta$  is a set that exhaustively collects all learnable parameters in the network. Since the joint distribution of  $\mathbf{x}$  and  $y$  is generally unknown in practice, the objective of expected risk minimization cannot be pursued directly and we opt to minimizing an empirical risk instead, using a set of  $N$  training samples  $\{(\mathbf{x}_i, y_i)\}$ , i.e.,  $l_{emp}(\Theta) := \sum_i l(f(\mathbf{x}_i; \Theta), y_i)/N$ . In spite of being a rule of thumb, there is always a gap between minimizing  $\tilde{l}(\Theta)$  and  $l_{emp}(\Theta)$ , and we can use the generalization error  $GE(g) = |\tilde{l}(\Theta) - l_{emp}(\Theta)|$  to characterize such a gap formally [14], [21].

It has been discovered that the margin and some other robustness measure of a classifier bound the generalization error of the classifier in theory [14], [22]. Naturally, we prefer machine learning models with lower generalization error

and it has been demonstrated that, for a  $(K, \epsilon(\mathcal{S}_n))$ -robust algorithm with  $l(f(\mathbf{x}, :), y) \leq M$  for all reasonable  $(\mathbf{x}, y)$ <sup>1</sup>, we have, with probability at least  $1 - \delta$ , it holds that

$$GE(g) \leq \epsilon(\mathcal{S}_n) + M \sqrt{\frac{2K \log(2) + 2 \log(1/\delta)}{N}}. \quad (2)$$

The result establishes connections between the generalization ability and robustness of classification models. Inspired the result, Sobolić et al. [23] proposed to bound the margin of a DNN such that both its robustness and generalization ability are guaranteed, and it was achieved by constraining the Jacobian matrix. For a classification model that was fed with an  $n$ -dimensional input each time and outputted  $k$  neurons before softmax, the  $n \times k$  Jacobian matrix should also be instance-specific and it was obtained by calculating the gradient of the function with respect to its input. Superior test-set accuracies were obtained using the Jacobian regularization. They also generalized the theoretical and empirical analyses to stable invariant classifiers (e.g., convolutional neural networks, CNNs) [24], [25], and methods that could enhance robustness to data variations are suggested. With all the facts, we know that, for Question 1, the large margin principle is beneficial to DNNs, in improving the generalization ability and robustness.

### 2.2 Large Margin Cannot be Trivially Obtained

Let us now turn to answering Question 2. Apparently, large margins cannot be naturally obtained with a normal training mechanism (i.e., using stochastic gradient descent to minimize just a cross-entropy loss) for DNN models in practice, otherwise their adversarial vulnerability would not be considered severe. Note that although it has been proved that linear networks trained on *separable data* using stochastic gradient descent converge to maximum margin solutions as  $t \rightarrow \infty$  [26], [27], [28], [29], it was also demonstrated that the convergence rate was extremely slow (e.g.,  $O(1/\log(t))$  [26] using the cross-entropy loss), making it hardly achievable in practice. Similar results can also be derived for *non-separable data* [30]. In addition to the results that were derived on the basis of the implicit bias of stochastic gradient descent, it has been shown that over parameterization also leads to improved margins [31]. Under an infinite network width regime, stochastic gradient descent of a two-layer network model leads to an inference function in a reproducing kernel Hilbert space of the neural tangent kernel [32], [33], [34], and it can be proved that a proper explicit regularization can indeed improve margins.

## 3 ACHIEVE LARGE MARGIN FOR DNNs

Now that we have answered the first two concerned questions in the previous section, we will attempt to answer the third question in this sections. Figure 2 is a summarization of what follows. Section 3.1 to 3.3 attempt to group training mechanisms that affect the margin with theoretical guarantees into several categories for better clarity.

1. Formal definition of  $(K, \epsilon(\mathcal{S}_n))$  can be found in [14], [23] and it characterizes how the training data is exploited by the algorithm.  $K$  is the number of sample partitions and  $\epsilon(\mathcal{S}_n)$  bounds the discrepancy between training and possible test losses in each partition.

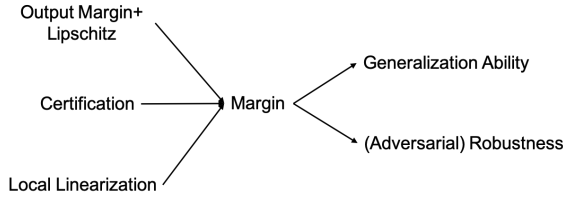


Fig. 2. The concerned classification margin is closely related to adversarial robustness and generalization ability, and we will introduce methods that aim to enlarge the margins for DNNs, by local linearization [35], [36], [37], certification [38], [39], [40], or relying on a margin in the decision space [41], [42], [43] in combination with some (possibly implicit) Lipschitz constraints.

### 3.1 Regularization on Lipschitz Constant and Output Margins

We know that conventional regularization strategies like weight decay [44] and dropout [45] do not help in enlarging the margin for nonlinear DNNs, which seems different from SVMs. It indicates that the hierarchical structure of DNNs is worthy further exploring, considering that the optimization of an SVM can be roughly considered as training the final layer of a DNN with weight decay.

Unlike for SVMs, derivation and calculation of a classification margin seem infeasible for nonlinear DNNs, and hence some approximations are pursued as surrogates. For instance, as mentioned, Sobolić et al. [23] took advantage of the fact that nonlinear DNNs were Lipschitz continuous and used the inequality

$$\|f(\mathbf{x}) - f(\mathbf{x}')\|_q \leq L_{p,q} \|\mathbf{x} - \mathbf{x}'\|_p \quad (3)$$

to constrain the margin, which was basically an  $l_p$  distance between two points in the input space. We know from Eq. (3) that the  $l_q$  distance between any two points in the prediction space is an essential factor of the classification margin. Fix  $\|f(\mathbf{x}) - f(\mathbf{x}')\|_q \leq \epsilon, \forall \mathbf{x}'$ , the classification margin can be guaranteed by minimizing a **Lipschitz constant** (i.e.,  $L_{p,q}$ ). It was also theoretically proven that the Lipschitz constant has direct relationships with both robustness [46], [47] and generalization [48], [49], [50] for a variety of nonlinear DNNs. Unsurprisingly, the Lipschitz constant and some “margins” in the model decision space, which laid the foundation of a spectrum of methods as will be introduced in the following paragraphs, jointly guarantee the (adversarial) robustness of DNNs. The norm of the Jacobian matrix, is then suggested to be penalized during training for improved robustness [23], [51], since it is actually a local Lipschitz constant fulfilling the inequality in Eq. (3) for any specific  $\mathbf{x}$  and it measures the sensitivity of a learning machine by definition [52]. In addition to the norm of the Jacobian matrix, the Euclidean norm of input gradient [53], the curvature of Hessian matrix [54], the spectral norm of weight matrices [55], and a cross-Lipschitz functional [46] can all be used to enhance the robustness or generalization ability of DNNs, and larger margins are simultaneously anticipated. In fact, these network properties are all closely related to the Lipschitz constant and have some chained inequality [51], [56]. In binary classification, some of the regularizations share the same essential ingredients (i.e., maximizing prediction confidence and minimizing local Lipschitz constants) in theory and show similar empirical results [56].

As mentioned, there is also a line of work focusing on “margins” in the decision or a representation space (i.e., **output margin**) of DNNs. For instance, Tang [57] proposed to replace the softmax cross-entropy loss with an SVM-derived one, leading to margin gain in the representation space characterized by the penultimate layer and a winning solution to the ICML’13 representation learning challenge [58]. Sun et al. [59] introduced margin-based penalties to the objective of training DNNs, motivated by theoretical analyses from the perspective of the margin bound. The triplet loss [60] that imposes a margin of representation distance between each positive sample pair and each negative sample pair was also considered, e.g., in FaceNet [61]. Since then, such output margins have been actively discussed in the face recognition community. Beside what was applied in FaceNet, angular softmax (A-Softmax) <sup>2</sup> and additive margin softmax (AM-Softmax) were developed and used in SphereFace [62] and CosineFace [42], respectively, to enhance the cross-entropy loss with novel softmax formulations. There are also ensemble soft-margin softmax (M-Softmax) [63], virtual softmax (V-Softmax) [64], large margin cosine loss (LMCL) [65], and additive angular margin loss (AAML) [43], just to name a few. They have achieved remarkable success in the task of face identification and verification, and some of them also showed promising accuracy for classifying natural scene images. The difference between these methods lie in the way of decomposing the cross-entropy loss. Table 1 summarizes and compares them. It is also worth mentioning that the cross-entropy loss itself can also be interpreted as a margin-based loss, with input-specific margins in the output space, and it was shown that enlarging such input-specific margins could be used as a regularization and led to improved test-set accuracy [66].

Similar ideas for enlarging output margins have also been considered in the task of speech recognition [67], [68], [69], in which computational modules like feedback connections [10] and self-attentions [70] often serve in the backbone models. It is expected to obtain deep feature representations that lead to the largest SVM margins. A two-stage pipeline was initially developed for training such models [67], [68], in which learnable parameters in the final and prior layers were updated separately. For more recent methods, automatic differentiation [71] was adopted, such that the feature representations and large margin SVM classifiers can be optimized jointly. Such a large margin principle was also utilized in few-shot learning [72], PU learning [73], [74], and anomaly detection [75].

While these methods considered classification margins of all training instances, it seems more efficient and reasonable to mainly focus on support vectors, just like in SVMs. In this spirit, Wang et al. [76] combined the margin-based softmax with hard example mining [77], [78], such that training focused more on harder samples which were considered more informative. Since only the “support vectors” were used for calculating gradients and updating parameters, the training process became more efficient. Being viewed as a functional abstraction of the training dataset, the set of “support vectors” has also been utilized to address catastrophic

<sup>2</sup>. See also large margin softmax (L-Softmax) in [41], which is very similar to A-Softmax.

TABLE 1

Large margin classification for face recognition. The methods differ from each other by introducing margins in angular, logit, or cosine spaces and by incorporating scaling factors (Mul) or additive terms (Add). Given  $W$  as a learnable matrix before softmax and  $h$  as the feature representation of a DNN, if  $\|W\|_2 = 1$  is fixed, then we can rewrite the linear transformation  $W^T h$  as  $\cos(\theta)\|h\|_2$  and incorporate a scaling factor  $\alpha$  and an additive term  $\beta$  into it to encourage the angular margin as  $\cos(\alpha\theta + \beta)\|h\|_2$ . Note that although M-Softmax targets to image classification, it is closely related to the other methods and thus we list it here as well.

Method	Test data	Margin	Add	Mul
L-Softmax [41]	image, face	Angular	✗	✓
A-Softmax [62]	face	Angular	✗	✓
M-Softmax [63]	image	Logit	✓	✗
V-Softmax [64]	image, face	Logit	✗	✗
AM-Softmax [42]	face	Cosine	✓	✗
LMCL [65]	face	Cosine	✓	✗
AAML [43]	face	Angular	✓	✗

forgetting [79] in incremental learning DNNs [80].

### 3.2 Local Linearization

The methods introduced in Section 3.1 enlarged “margins” mostly in the decision space of DNN models. Efforts might also be devoted in other representation spaces [81], however, as discussed [82], owing to the distance distortions between input and representation spaces, the classification margins in the input space of DNNs were not necessarily maximized by methods in this category<sup>3</sup>. An et al. [82] hence further enforced the transformations of middle layers of a DNN to be contraction mappings, in order to achieve large classification margins. One step further, Bansal et al. [85] advocated a similar learning objective to that of SVMs, in which the method of Lagrange multipliers is utilized. Both of them can be seen as using layer-wise approximations to restricting some whole network property (e.g., the Lipschitz constant). Recently, independent work from Yan et al. [35] and Elsayed et al. [37] proposed to enlarge classification margins via **local linearization** and reasonable approximations. In fact, by simply rewriting the constraint  $f(\mathbf{x})f(\mathbf{x} + \mathbf{z}) < 0$  using Taylor’s approximation of  $f$  with respect to  $\mathbf{x}$ , one can obtain a closed-form solution to the worst-case perturbations even for nonlinear DNNs. Elsayed et al. proposed to maximize

$$\frac{|f(\mathbf{x})|}{\|\nabla_{\mathbf{x}} f(\mathbf{x})\|_q}, \quad (4)$$

in which  $\|\cdot\|_q$  was the dual norm of  $\|\cdot\|_p$ , with  $1/p + 1/q = 1$ . In essence, the method can be considered as a simplified and efficient version of Yan et al.’s [35], while Yan et al.’s method took one step further and followed a similar mechanism to DeepFool [18]. Specifically, an iterative local linearization was utilized by DeepFool and Yan et al.’s method to pursue approximations to evaluating the classification margins [35], which can be more accurate yet more computational demanding. Shortly after, Ding et al. [86] also discussed large

3. Although it has been theoretically shown that enlarging the output margin is beneficial to the generalization ability of DNNs, along with constrained classifier norms [48], [83] or constrained complexity of each layer [84]

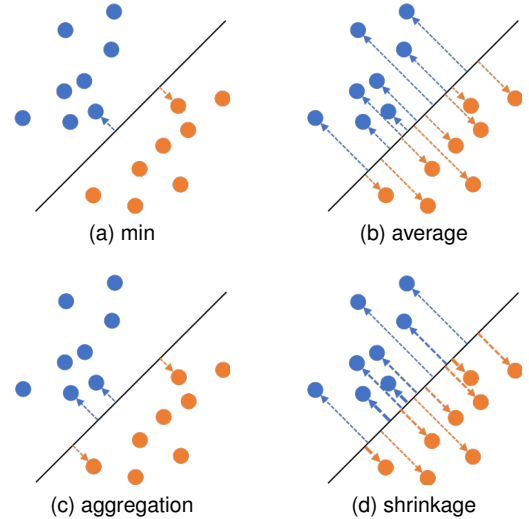


Fig. 3. Different choices of the key components result in different penalty on instance-specific margins. More specifically, (a) min considers only the worst case in the training set, (b)  $\sum$  combined with an identity function treat all training samples equally, (c) and (d) incorporates the aggregation function and shrinkage function, respectively.

margin training for DNNs, and they proposed to incorporate

$$\sum_{i \in S^+} \max\{0, m_{\max} - \hat{m}_{\mathbf{x}_i, f}\} \quad (5)$$

into the learning objective of DNNs, in which  $\hat{m}_{\mathbf{x}_i, f}$  was an estimation of the instance-specific margin  $m_{\mathbf{x}_i, f}$  for  $\mathbf{x}_i$ ,  $S^+$  indicated the set of correctly classified training samples using the network model, and  $m_{\max} > 0$  was a hyperparameter. To well approximate  $m_{\mathbf{x}, f}$ , they took advantage of the PGD attack [87] and adopted its variant as a proxy of the “shortest successful perturbation”.

#### 3.2.1 Key Components and More Discussions

In existing work that explicitly incorporates margin-based regularizers into learning objectives, given the set of training samples  $\{(\mathbf{x}_i, y_i)\}$ , one might write the regularization term as

$$\frac{1}{|\mathcal{T}|} \sum_{i \in \mathcal{T}} r(-m_{\mathbf{x}_i, f}) \quad \text{or} \quad \max_{i \in \mathcal{T}} r(-m_{\mathbf{x}_i, f}), \quad (6)$$

in which  $r(\cdot)$  is a monotonically increasing function. It has been discussed that the specific choice of  $\max$  (that focuses on the minimum margin solely) and  $\sum$  indicates preference in improving robustness or generalization ability [36], [88]. The function  $r(\cdot)$  is introduced to optionally put more stress on samples with smaller margins.  $\mathcal{T}$  can be used to enlarge margins for correctly classified training samples only, since the margin for incorrectly classified samples are not well defined. They form the three key components in developing different margin-based regularizers.

Somewhat surprisingly, Wu et al. [88] showed that there might exist an interesting trade-off between minimum margins and average margins. It is uncontroversial, at least for SVMs, that maximizing the minimum margin (i.e., choosing the formulation with  $\max$  in Eq. (6)) leads to improved generalization ability, while the average-based term in Eq. (6)

favors adversarial robustness<sup>4</sup>. In fact, we typically use the average magnitude to evaluate the (adversarial) robustness of DNNs [18], [90], [91], [92]. However, it is also discovered that using an identity function for  $r(\cdot)$  as  $-\sum_{i \in \mathcal{T}} m_{\mathbf{x}_i, f} / |\mathcal{T}|$  results in poor prediction accuracy on benign inputs. The regularizer is found to be dominated by some “extremely robust” samples so that if all the samples are treated equally in the regularizer, it would be difficult to well trade off the benign-set accuracy and adversarial robustness. A nonlinear “shrinkage” function that penalize more on samples with small margins can be introduced to relieve the problem [35], [93].

Although using min or max instead of the average operation in the regularizer aligns better with the generalization ability, it may lead to slower training convergence since at most one (instance-specific) margin in a whole batch of samples is effectively optimized at each training iteration. Yan et al. [36] proposed to address this issue by incorporating an “aggregation” function that aggregates some samples in the batch rather than using only one.

For  $\mathcal{T}$ , some methods chose it to be a set of correctly classified samples [35], [36], [86] while the others simply used the whole training set. For datasets where nearly 100% training accuracy can be obtained, the two options actually lead to similar performance. See Figure 3 for a comparison of different settings in the functions.

### 3.3 DNN Certification

A related and surging category of methods for DNN **certification** (also known as DNN verification) was also widely explored in the community [38], [39], [40], [94], [95], [96], [97], [98], [99], [100]. These DNN certification methods aimed at maximizing the volume of the hyper-sphere centered at each training instance, in which all data points are predicted into the same class. In fact, the radius of such a certified hyper-sphere is actually a lower bound of the classification margin, therefore the techniques could also be regarded as opting to encouraging large classification margins.

The certifications of DNNs are normally rigorous such that the margins and DNN robustness can be theoretically guaranteed. One of their demerits might be the high computational complexity. In fact, it is challenging for most of them to be generalized to large networks on large datasets (e.g., ImageNet [101]). For fast certification, Lee et al. [102] and Croce et al. [103] proposed to expand linear regions where training samples reside. The linear regions can be smaller than the certified hyper-spheres, and they are related to the classification margins similarly. The relationship between margins and the Lipschitz constant is also exploited, to achieve better certification efficiency [104].

## 4 DATA AUGMENTATION AND DNN COMPRESSION MAY AFFECT MARGINS

In addition to the methods introduced in Section 3, there exist other methods that possibly achieve large classification margins as some sort of a byproduct. In general, methods

4. It has also been shown that in addition to the average margin which characterizes the first-order statistics of the margin distribution, the variance of the margin is also of importance [89]

that benefit the generalization ability and test accuracy of DNNs may unintentionally enlarge margins to some extent. From this point of view, it has been discussed under what condition(s) can data augmentation lead to margin improvement [105]. An achieved result is that, for linear classifiers or linearly separable data, polynomially many more samples are required for a very specific data augmentation strategy to obtain optimal margins. Adversarial training [87] is often also regarded as a data augmentation strategy, and we know that it has to enlarge margin to guarantee performance. In fact, it has been shown that adversarial training converges to maximum margin solutions faster than normal training [106]. Other augmentations include cutout [107] and mixup [108] may also be related with enlarged margins, considering their success in improving the generalization ability of DNNs, and we encourage future work to explore along this direction.

Though it lacks obvious evidence that can demonstrate the relationship between DNN margins and other learning technologies, we would like to discuss directions that can possibly be explored. The first set of technologies that attract our attention is DNN compression, including network pruning [109], [110], [111] and quantization [112], [113], since improved or at least similar test-set accuracy can be achieved with significantly fewer learnable parameters (in bit for quantization) using these methods [113], [114]. It is also discovered that such network compression leads to improved adversarial robustness in certain circumstances [92], [115]. As have been introduced, both the generalization ability and robustness are closely related to the classification margins, thus we conjecture that it is possible that compression along with re-training also help to achieve models with enlarged margins.

## 5 ESTIMATION OF THE MARGIN

Classification margins have been used for a variety of goals, e.g., improving and estimating DNN model robustness and generalization ability [35], [36], [37], [86], [116]. However, it is still an open problem to find an accurate approximation for the DNN classification margin. As mentioned, the radius of a certified hyper-sphere bounds the margin from below, while the adversarial examples act as upper bounds of the margin. Hence, taking advantage of DNN certification [38], [39], [40], [94], [95], [96], [97], [98], [99], [100] and adversarial attacks [87], [91], one can reasonably estimate the range where the classification margin resides in. Other methods for estimating the margin, probably not rigorously, can also be found in Section 3.

## 6 CONCLUSION

In this paper, we have surveyed recent research efforts on classification margin for (nonlinear) DNNs. Unlike for SVMs, the studies are more challenging for DNNs on account of their hierarchical structure and complex nonlinearity. We have revisited some work in the last century and highlight the focus of this paper in the first section, and we have then summarized connections between the margin, generalization, and robustness, mostly from a theoretical point of view, which highlights the importance of large

margin even in the state-of-the-art DNN models. We have reviewed methods that target at large margin DNNs over the past few years, and we have categorized them into groups, in a comprehensive but summarized manner. We managed to shorten the paper such that crucial spirit of large margin learning and related methods could be better emphasized. We have shared our view on the key components of current winning methods and point to directions that can possibly be explored.

## ACKNOWLEDGMENTS

This work is funded by the National Key Research and Development Program of China (No. 2018AAA0100701) and a grant from the Guoqiang Institute, Tsinghua University.

## REFERENCES

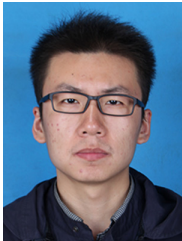
- [1] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [2] V. Vapnik, "An overview of statistical learning theory," *IEEE transactions on neural networks*, vol. 10, no. 5, pp. 988–999, 1999.
- [3] —, *The nature of statistical learning theory*. Springer science & business media, 2013.
- [4] K. Crammer and Y. Singer, "On the algorithmic implementation of multiclass kernel-based vector machines," *Journal of machine learning research*, vol. 2, no. Dec, pp. 265–292, 2001.
- [5] I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun, "Large margin methods for structured and interdependent output variables," *Journal of machine learning research*, vol. 6, no. Sep, pp. 1453–1484, 2005.
- [6] J. Zhu, S. Rosset, R. Tibshirani, and T. J. Hastie, "1-norm support vector machines," in *Advances in neural information processing systems*, 2004, pp. 49–56.
- [7] A. Cotter, S. Shalev-Shwartz, and N. Srebro, "Learning optimally sparse support vector machines," in *International Conference on Machine Learning*, 2013, pp. 266–274.
- [8] A. J. Smola, P. J. Bartlett, D. Schuurmans, and B. Schölkopf, *Advances in large margin classifiers*. MIT press, 2000.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [10] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *ICASSP*, 2013.
- [11] J. Shawe-Taylor, N. Cristianini *et al.*, *Kernel methods for pattern analysis*. Cambridge university press, 2004.
- [12] P. Bartlett and J. Shawe-Taylor, "Generalization performance of support vector machines and other pattern classifiers," *Advances in Kernel methods - support vector learning*, pp. 43–54, 1999.
- [13] B. Schölkopf, A. J. Smola, R. C. Williamson, and P. L. Bartlett, "New support vector algorithms," *Neural computation*, vol. 12, no. 5, pp. 1207–1245, 2000.
- [14] H. Xu, C. Caramanis, and S. Mannor, "Robustness and regularization of support vector machines," *Journal of Machine Learning Research*, vol. 10, no. Jul, pp. 1485–1510, 2009.
- [15] H. Xu and S. Mannor, "Robustness and generalization," *Machine learning*, vol. 86, no. 3, pp. 391–423, 2012.
- [16] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *ICLR*, 2014.
- [17] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *ICLR*, 2015.
- [18] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "DeepFool: a simple and accurate method to fool deep neural networks," in *CVPR*, 2016.
- [19] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in *ICLR*, 2017.
- [20] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille, "Adversarial examples for semantic segmentation and object detection," in *ICCV*, 2017.
- [21] K. Kawaguchi, L. P. Kaelbling, and Y. Bengio, "Generalization in deep learning," *arXiv preprint arXiv:1710.05468*, 2017.
- [22] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *Advances in Neural Information Processing Systems*, 2019, pp. 1567–1578.
- [23] J. Sokolić, R. Giryes, G. Sapiro, and M. R. Rodrigues, "Robust large margin deep neural networks," *IEEE Transactions on Signal Processing*, vol. 65, no. 16, pp. 4265–4280, 2017.
- [24] J. Sokolic, R. Giryes, G. Sapiro, and M. Rodrigues, "Generalization error of invariant classifiers," in *AISTATS*, 2017, pp. 1094–1103.
- [25] J. Huang, Q. Qiu, G. Sapiro, and R. Calderbank, "Discriminative robust transformation learning," in *NeurIPS*, 2015, pp. 1333–1341.
- [26] D. Soudry, E. Hoffer, M. S. Nacson, S. Gunasekar, and N. Srebro, "The implicit bias of gradient descent on separable data," *Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2822–2878, 2018.
- [27] S. Gunasekar, J. D. Lee, D. Soudry, and N. Srebro, "Implicit bias of gradient descent on linear convolutional networks," in *NeurIPS*, 2018, pp. 9461–9471.
- [28] Z. Ji and M. Telgarsky, "Gradient descent aligns the layers of deep linear networks," in *ICLR*, 2018.
- [29] M. S. Nacson, J. Lee, S. Gunasekar, P. H. P. Savarese, N. Srebro, and D. Soudry, "Convergence of gradient descent on separable data," in *AISTATS*, 2019, pp. 3420–3428.
- [30] Z. Ji and M. Telgarsky, "The implicit bias of gradient descent on nonseparable data," in *Conference on Learning Theory*, 2019, pp. 1772–1798.
- [31] C. Wei, J. D. Lee, Q. Liu, and T. Ma, "Regularization matters: Generalization and optimization of neural nets v.s. their induced kernel," in *NeurIPS*, 2019.
- [32] A. Daniely, "Sgd learns the conjugate kernel class of the network," in *Advances in Neural Information Processing Systems*, 2017, pp. 2422–2430.
- [33] A. Jacot, F. Gabriel, and C. Hongler, "Neural tangent kernel: Convergence and generalization in neural networks," in *Advances in neural information processing systems*, 2018, pp. 8571–8580.
- [34] S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, and R. Wang, "On exact computation with an infinitely wide neural net," in *Advances in Neural Information Processing Systems*, 2019, pp. 8141–8150.
- [35] Z. Yan, Y. Guo, and C. Zhang, "Deep defense: Training dnns with improved adversarial robustness," in *NeurIPS*, 2018.
- [36] —, "Adversarial margin maximization networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [37] G. F. Elsayed, D. Krishnan, H. Mobahi, K. Regan, and S. Bengio, "Large margin deep networks for classification," in *NeurIPS*, 2018.
- [38] V. Tjeng, K. Xiao, and R. Tedrake, "Evaluating robustness of neural networks with mixed integer programming," *arXiv preprint arXiv:1711.07356*, 2017.
- [39] E. Wong and Z. Kolter, "Provable defenses against adversarial examples via the convex outer adversarial polytope," in *ICML*, 2018, pp. 5283–5292.
- [40] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *ICML*, 2019, pp. 1310–1320.
- [41] W. Liu, Y. Wen, Z. Yu, and M. Yang, "Large-margin softmax loss for convolutional neural networks," in *ICML*, 2016, pp. 507–516.
- [42] F. Wang, J. Cheng, W. Liu, and H. Liu, "Additive margin softmax for face verification," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.
- [43] J. Deng, J. Guo, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *CVPR*, 2019.
- [44] A. Krogh and J. A. Hertz, "A simple weight decay can improve generalization," in *NeurIPS*, 1992, pp. 950–957.
- [45] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.
- [46] M. Hein and M. Andriushchenko, "Formal guarantees on the robustness of a classifier against adversarial manipulation," in *NeurIPS*, 2017, pp. 2266–2276.
- [47] T.-W. Weng, H. Zhang, P.-Y. Chen, J. Yi, D. Su, Y. Gao, C.-J. Hsieh, and L. Daniel, "Evaluating the robustness of neural networks: An extreme value theory approach," in *ICLR*, 2018.
- [48] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky, "Spectrally-normalized margin bounds for neural networks," in *NeurIPS*, 2017, pp. 6240–6249.

- [49] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," in *NeurIPS*, 2017, pp. 5947–5956.
- [50] C. Wei and T. Ma, "Data-dependent sample complexity of deep neural networks via lipschitz augmentation," in *NeurIPS*, 2019.
- [51] D. Jakubovitz and R. Giryes, "Improving dnn robustness to adversarial attacks using jacobian regularization," in *ECCV*, 2018.
- [52] R. Novak, Y. Bahri, D. A. Abolafia, J. Pennington, and J. Sohl-Dickstein, "Sensitivity and generalization in neural networks: an empirical study," in *ICLR*, 2018.
- [53] A. S. Ross and F. Doshi-Velez, "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients," in *AAAI*, 2018.
- [54] S.-M. Moosavi-Dezfooli, A. Fawzi, J. Uesato, and P. Frossard, "Robustness via curvature regularization, and vice versa," in *CVPR*, 2019.
- [55] Y. Yoshida and T. Miyato, "Spectral norm regularization for improving the generalizability of deep learning," *arXiv preprint arXiv:1705.10941*, 2017.
- [56] Y. Guo, L. Chen, Y. Chen, and C. Zhang, "On connections between regularizations for improving dnn robustness," *IEEE transactions on pattern analysis and machine intelligence*, 2020.
- [57] Y. Tang, "Deep learning using linear support vector machines," *ICML Workshop*, 2013.
- [58] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville, M. Mirza, B. Hamner, W. Cukierski, Y. Tang, D. Thaler, D.-H. Lee *et al.*, "Challenges in representation learning: A report on three machine learning contests," *arXiv preprint arXiv:1307.0414*, 2013.
- [59] S. Sun, W. Chen, L. Wang, X. Liu, and T.-Y. Liu, "On the depth of deep neural networks: A theoretical view." in *AAAI*, 2016, pp. 2066–2072.
- [60] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *Journal of Machine Learning Research*, vol. 10, no. Feb, pp. 207–244, 2009.
- [61] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *CVPR*, 2015, pp. 815–823.
- [62] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *CVPR*, 2017.
- [63] X. Wang, S. Zhang, Z. Lei, S. Liu, X. Guo, and S. Z. Li, "Ensemble soft-margin softmax loss for image classification," in *IJCAI*, 2018.
- [64] B. Chen, W. Deng, and H. Shen, "Virtual class enhanced discriminative embedding learning," in *Advances in Neural Information Processing Systems*, 2018, pp. 1942–1952.
- [65] H. Wang, Y. Wang, Z. Zhou, X. Ji, Z. Li, D. Gong, J. Zhou, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *CVPR*, 2018, pp. 5265–5274.
- [66] T. Kobayashi, "Large margin in softmax cross-entropy loss." in *BMVC*, 2019, p. 139.
- [67] S.-X. Zhang, C. Liu, K. Yao, and Y. Gong, "Deep neural support vector machines for speech recognition," in *ICASSP*, 2015, pp. 4275–4279.
- [68] S.-X. Zhang, R. Zhao, C. Liu, J. Li, and Y. Gong, "Recurrent support vector machines for speech recognition," in *ICASSP*, 2016, pp. 5885–5889.
- [69] P. Wang, J. Cui, C. Weng, and D. Yu, "Large margin training for attention based end-to-end speech recognition," in *Interspeech*, 2019, pp. 246–250.
- [70] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [71] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *NIPS Autodiff Workshop*, 2017.
- [72] Y. Wang, X.-M. Wu, Q. Li, J. Gu, W. Xiang, L. Zhang, and V. O. Li, "Large margin few-shot learning," *arXiv preprint arXiv:1807.02872*, 2018.
- [73] C. Elkan and K. Noto, "Learning classifiers from only positive and unlabeled data," in *SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2008, pp. 213–220.
- [74] T. Gong, G. Wang, J. Ye, Z. Xu, and M. Lin, "Margin based pu learning," in *AAAI*, 2018.
- [75] W. Liu, W. Luo, Z. Li, P. Zhao, and S. Gao, "Margin learning embedded prediction for video anomaly detection with a few anomalies," in *IJCAI*, 2019, pp. 3023–3030.
- [76] X. Wang, S. Wang, S. Zhang, T. Fu, H. Shi, and T. Mei, "Support vector guided softmax loss for face recognition," *arXiv preprint arXiv:1812.11317*, 2018.
- [77] A. Shrivastava, A. Gupta, and R. Girshick, "Training region-based object detectors with online hard example mining," in *CVPR*, 2016, pp. 761–769.
- [78] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017, pp. 2980–2988.
- [79] R. Kemker, M. McClure, A. Abitino, T. L. Hayes, and C. Kanan, "Measuring catastrophic forgetting in neural networks," in *AAAI*, 2018.
- [80] Y. Li, Z. Li, L. Ding, Y. Pan, C. Huang, Y. Hu, W. Chen, and X. Gao, "Supportnet: solving catastrophic forgetting in class incremental learning with support data," *arXiv preprint arXiv:1806.02942*, 2018.
- [81] Y. Zhong and W. Deng, "Adversarial learning with margin-based triplet embedding regularization," *arXiv preprint arXiv:1909.09481*, 2019.
- [82] S. An, M. Hayat, S. H. Khan, M. Bennamoun, F. Boussaid, and F. Sohel, "Contractive rectifier networks for nonlinear maximum margin classification," in *ICCV*, 2015, pp. 2515–2523.
- [83] B. Neyshabur, S. Bhojanapalli, and N. Srebro, "A pac-bayesian approach to spectrally-normalized margin bounds for neural networks," *arXiv preprint arXiv:1707.09564*, 2017.
- [84] C. Wei and T. Ma, "Improved sample complexities for deep networks and robust classification via an all-layer margin," in *ICML*, 2020.
- [85] Y. Bansal, M. Advani, D. Cox, and A. Saxe, "Minnorm training: an algorithm for training overcomplete deep neural networks," *arXiv preprint arXiv:1806.00730*, 2018.
- [86] G. W. Ding, Y. Sharma, K. Y. C. Lui, and R. Huang, "Max-margin adversarial (mma) training: Direct input space margin maximization through adversarial training," *arXiv preprint arXiv:1812.02637*, 2019.
- [87] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *ICLR*, 2018.
- [88] K. Wu and Y. Yu, "Understanding adversarial robustness: The trade-off between minimum and average margin," *arXiv preprint arXiv:1907.11780*, 2019.
- [89] T. Zhang and Z.-H. Zhou, "Optimal margin distribution machine," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 6, pp. 1143–1156, 2019.
- [90] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *CVPR*, 2017.
- [91] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.
- [92] Y. Guo, C. Zhang, C. Zhang, and Y. Chen, "Sparse dnns with improved adversarial robustness," in *NeurIPS*, 2018, pp. 242–251.
- [93] D. Stutz, M. Hein, and B. Schiele, "Confidence-calibrated adversarial training: Towards robust models generalizing beyond the attack used during training," *arXiv preprint arXiv:1910.06259*, 2019.
- [94] C.-H. Cheng, G. Nührenberg, and H. Ruess, "Maximum resilience of artificial neural networks," in *International Symposium on Automated Technology for Verification and Analysis*, 2017, pp. 251–268.
- [95] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer, "Reluplex: An efficient smt solver for verifying deep neural networks," in *International Conference on Computer Aided Verification*, 2017, pp. 97–117.
- [96] K. Dvijotham, R. Stanforth, S. Gowal, T. A. Mann, and P. Kohli, "A dual approach to scalable verification of deep networks." in *UAI*, 2018, pp. 550–559.
- [97] G. Singh, T. Gehr, M. Mirman, M. Püschel, and M. Vechev, "Fast and effective robustness certification," in *NeurIPS*, 2018, pp. 10 802–10 813.
- [98] A. Boopathy, T.-W. Weng, P.-Y. Chen, S. Liu, and L. Daniel, "Cnn-cert: An efficient framework for certifying robustness of convolutional neural networks," in *AAAI*, vol. 33, 2019, pp. 3240–3247.
- [99] H. Zhang, T.-W. Weng, P.-Y. Chen, C.-J. Hsieh, and L. Daniel, "Efficient neural network robustness certification with general activation functions," in *NeurIPS*, 2018, pp. 4939–4948.

- [100] A. Fromherz, K. Leino, M. Fredrikson, B. Parno, and C. Păsăreanu, "Fast geometric projections for local robustness certification," in *ICLR*, 2021.
- [101] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "Imagenet large scale visual recognition challenge," *IJCV*, 2015.
- [102] G.-H. Lee, D. Alvarez-Melis, and T. S. Jaakkola, "Towards robust, locally linear deep networks," in *ICLR*, 2019.
- [103] F. Croce, M. Andriushchenko, and M. Hein, "Provable robustness of relu networks via maximization of linear regions," in *AISTATS*, 2019.
- [104] Y. Tsuzuku, I. Sato, and M. Sugiyama, "Lipschitz-margin training: Scalable certification of perturbation invariance for deep neural networks," in *NeurIPS*, 2018, pp. 6541–6550.
- [105] S. Rajput, Z. Feng, Z. Charles, P.-L. Loh, and D. Papailiopoulos, "Does data augmentation lead to positive margin?" in *ICML*, 2019.
- [106] Z. Charles, S. Rajput, S. Wright, and D. Papailiopoulos, "Convergence and margin of adversarial training on separable data," *arXiv preprint arXiv:1905.09209*, 2019.
- [107] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *arXiv preprint arXiv:1708.04552*, 2017.
- [108] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [109] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *NeurIPS*, 2015, pp. 1135–1143.
- [110] Y. Guo, A. Yao, and Y. Chen, "Dynamic network surgery for efficient dnns," in *NeurIPS*, 2016, pp. 1379–1387.
- [111] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning structured sparsity in deep neural networks," in *NeurIPS*, 2016, pp. 2074–2082.
- [112] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," *arXiv preprint arXiv:1602.02830*, 2016.
- [113] A. Zhou, A. Yao, Y. Guo, L. Xu, and Y. Chen, "Incremental network quantization: Towards lossless cnns with low-precision weights," in *ICLR*, 2017.
- [114] S. Han, J. Pool, S. Narang, H. Mao, E. Gong, S. Tang, E. Elsen, P. Vajda, M. Paluri, J. Tran *et al.*, "Dsd: Dense-sparse-dense training for deep neural networks," in *ICLR*, 2016.
- [115] A. Galloway, G. W. Taylor, and M. Moussa, "Attacking binarized neural networks," in *ICLR*, 2017.
- [116] Y. Jiang, D. Krishnan, H. Mobahi, and S. Bengio, "Predicting the generalization gap in deep networks with margin distributions," in *ICLR*, 2019.



**Changshui Zhang** received the B.E. degree in mathematics from Peking University, Beijing, China, in 1986, and the M.S. and Ph.D. degrees in control science and engineering from Tsinghua University, Beijing, in 1989 and 1992, respectively. In 1992, he joined the Department of Automation, Tsinghua University, where he is currently a professor. His research interests include pattern recognition and machine learning. He is a Fellow member of the IEEE.



**Yiwen Guo** received the B.E. degree from Wuhan University, Wuhan, China, in 2011, and the Ph.D. degree from Tsinghua University, Beijing, China in 2016. He is currently a research scientist at ByteDance AI Lab. Prior to this, he was a staff research scientist at Intel Labs. His current research interests include computer vision, machine learning, and security.