

# Learning Through Margins

Guangzeng Xie 1901111419

April 20, 2020

# Contents

Large Margin Deep Networks for Classification

Predicting the Generalization Gap in Deep Networks with Margin Distributions

Improved Sample Complexities for Deep Networks and Robust Classification via an All-Layer Margin

# Large Margin Deep Networks for Classification

- ▶ Consider a classification problem with  $n$  classes.
- ▶ Suppose we use a function  $f_i : \mathcal{X} \rightarrow \mathbb{R}$ , for  $i = 1, \dots, n$  that generates a prediction score for classifying the input vector  $\mathbf{x} \in \mathcal{X}$  to class  $i$ .
- ▶ The normal margin on example  $(\mathbf{x}, y)$  is defined by  $\gamma(\mathbf{f}, \mathbf{x}, y) = \max\{0, f_y(\mathbf{x}) - \max_{i \neq y} f_i(\mathbf{x})\}$ .
- ▶ This paper defined a novel loss to penalize the displacement of each  $\mathbf{x}$  to satisfy the margin constraint for separating class  $y$  from class  $i$  ( $i \neq y$ ).

- ▶ This paper proposed following optimization problem:

$$\min \sum_k \mathcal{A}_{i \neq y_k} (\max\{0, \gamma + d(\mathbf{f}, \mathbf{x}_k, (i, y_k)) \operatorname{sign}(f_i(\mathbf{x}_k) - f_{y_k}(\mathbf{x}_k))\}),$$

where  $\gamma$  is a hyper-parameter,  $\mathcal{A}$  is an aggregation operator like  $\sum$  and  $\max$ .

- ▶  $d(\mathbf{f}, \mathbf{x}, (i, j))$  is the distance of a point  $\mathbf{x}$  to the decision boundary  $\{\mathbf{x} : f_i(\mathbf{x}) = f_j(\mathbf{x})\}$  and is defined as

$$d(\mathbf{f}, \mathbf{x}, (i, j)) \triangleq \min_{\delta} \|\delta\|_p \quad \text{s.t.} \quad f_i(\mathbf{x} + \delta) = f_j(\mathbf{x} + \delta).$$

- The authors used a linearization technique to estimate  $d$  by

$$\bar{d}(\mathbf{f}, \mathbf{x}, (i, j)) \triangleq \min_{\boldsymbol{\delta}} \|\boldsymbol{\delta}\|_p \quad \text{s.t.} \quad f_i(\mathbf{x}) + \langle \boldsymbol{\delta}, \nabla_{\mathbf{x}} f_i(\mathbf{x}) \rangle = f_j(\mathbf{x}) + \langle \boldsymbol{\delta}, \nabla_{\mathbf{x}} f_j(\mathbf{x}) \rangle,$$

that is

$$\bar{d}(\mathbf{f}, \mathbf{x}, (i, j)) = \frac{|f_i(\mathbf{x}) - f_j(\mathbf{x})|}{\|\nabla_{\mathbf{x}} f_i(\mathbf{x}) - \nabla_{\mathbf{x}} f_j(\mathbf{x})\|_q},$$

where  $1/p + 1/q = 1$ .

- ▶ In fact, the activations at each intermediate layer could be interpreted as some intermediate representation of the data for the following part of the network. Thus, we can define the margin based on any intermediate representation and the ultimate decision boundary.
- ▶ Finally, the optimization problem for the deep networks can be written as:

$$\min \sum_{l,k} \mathcal{A}_{i \neq y_k} \left( \max \left\{ 0, \gamma_l + \frac{f_l(\mathbf{x}_k) - f_{y_k}(\mathbf{x}_k)}{\varepsilon + \|\nabla_{\mathbf{h}_l} f_l(\mathbf{x}_k) - \nabla_{\mathbf{h}_l} f_{y_k}(\mathbf{x}_k)\|_q} \right\} \right).$$

# Experiments

In experiments, to reduce the computational cost, the authors

- ▶ choose a subset of the total number of classes,
- ▶ treat the denominator as a constant,
- ▶ choose 5 evenly spaced layers (input layer, output layer and 3 other convolutional layers in the middle) across the network.

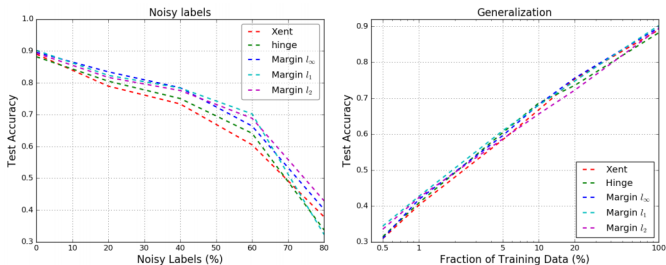


Figure 4: Performance of CIFAR-10 models on noisy data (left) and limited data (right).

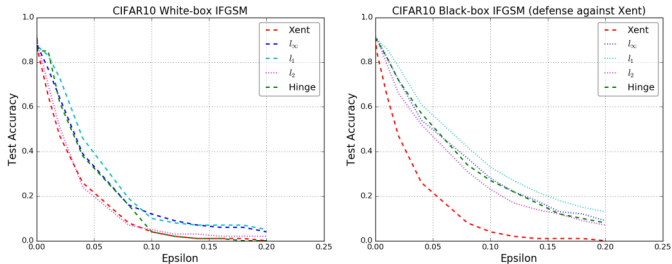
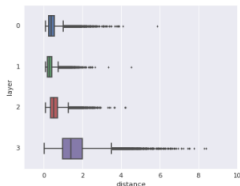
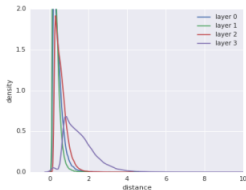


Figure 5: Performance of CIFAR-10 models on IFGSM adversarial examples.

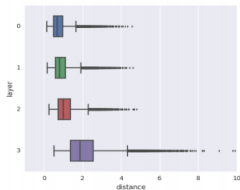
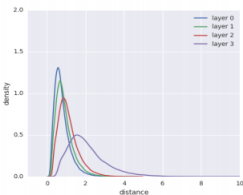


# Predicting the Generalization Gap in Deep Networks with Margin Distributions

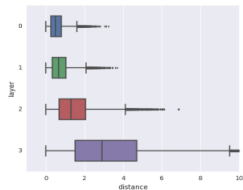
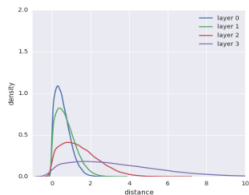
Test Acc.: 55.2%



Test Acc.: 70.6%



Test Acc.: 85.1%



- ▶ The normalized margin is specified by

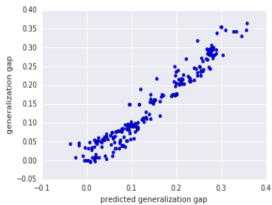
$$\hat{d}(\mathbf{f}, \mathbf{x}_k, (i, j)) = \frac{d(\mathbf{f}, \mathbf{x}_k, (i, j))}{\sqrt{\nu(\mathbf{x})}},$$

where  $\nu(\mathbf{x}) = \text{Tr}(\frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^\top)$ .

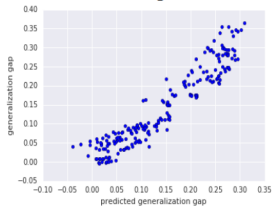
- ▶ The signatures of the distribution of normalized margins used in paper are
  - ▶ first five moments, or
  - ▶ the median  $Q_2$ , first quartile  $Q_1$  and third quartile  $Q_3$ , the upper fence  $\max\{\hat{d}_m : \hat{d}_m \leq Q_3 + 1.5\text{IQR}\}$  and the lower fence  $\min\{\hat{d}_m : \hat{d}_m \leq Q_1 - 1.5\text{IQR}\}$  where  $\text{IQR} = Q_3 - Q_1$ .

- ▶ Total signature  $\theta$  used in the paper to predict the generalization gap is a vector with dimension 20 (four evenly-spaced layers (input, and 3 hidden layers) with 5 signatures).
- ▶ The prediction model is a linear model:  $\hat{g} = \mathbf{a}^\top \phi(\theta) + b$ , where  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is a function applied element-wise to  $\theta$ .

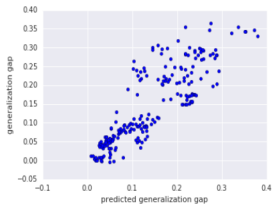
Norm. Margin 20D



Norm. Margin 4D



Bartlett Margin 5D



# Improved Sample Complexities for Deep Networks and Robust Classification via an All-Layer Margin

- ▶ Suppose that the classifier  $F(\mathbf{x}) = f_k \circ \dots \circ f_1(\mathbf{x})$  is computed by composing  $k$  functions  $f_k, \dots, f_1$ .
- ▶ Let  $\delta_k, \dots, \delta_1$  denote perturbations intended to be applied at each hidden layer. We recursively define the perturbed network output  $F(\mathbf{x}, \delta_1, \dots, \delta_k)$  by

$$h_1(\mathbf{x}, \delta) = f_1(\mathbf{x}) + \|\mathbf{x}\|_2 \delta_1,$$

$$h_i(\mathbf{x}, \delta) = f_i(h_{i-1}(\mathbf{x}, \delta)) + \|h_{i-1}(\mathbf{x}, \delta)\|_2 \delta_i,$$

$$F(\mathbf{x}, \delta) = h_k(\mathbf{x}, \delta).$$

- ▶ The all-layer margin will now be defined as the minimum norm of required to make the classifier misclassify the input:

$$m_F(\mathbf{x}, y) \triangleq \min_{\delta} \sqrt{\sum_{i=1}^k \|\delta_i\|_2^2} \quad \text{s.t.} \quad \arg \max_j F(\mathbf{x}, \delta)_j \neq y.$$

# Theoretical results

**Theorem 2.1** (Simplified version of Theorem A.1). *In the above setting, with probability  $1 - \delta$  over the draw of the training data, all classifiers  $F \in \mathcal{F}$  which achieve training error 0 satisfy*

$$\mathbb{E}_P[\ell_{0-1}(F(x), y)] \lesssim \frac{\sum_i \mathcal{C}_i}{\sqrt{n}} \sqrt{\mathbb{E}_{(x,y) \sim P_n} \left[ \frac{1}{m_F(x,y)^2} \right]} \log^2 n + \zeta$$

where  $\zeta \triangleq O\left(\frac{\log(1/\delta) + \log n}{n}\right)$  is a low-order term.

# Empirical Application of the All-Layer Margin

We can consider another objective:

$$\min_{\Theta} \max_{\delta} G(\delta, \Theta; \mathbf{x}, y) \triangleq \ell(F_{\Theta}(\mathbf{x}, \delta), y) - \lambda \|\delta\|_2^2$$

---

**Algorithm 1** All-layer Margin Optimization (AMO)

---

**procedure** PERTURBEDUPDATE(minibatch  $B = \{(x_i, y_i)\}_{i=1}^b$ , current parameters  $\Theta$ )

Initialize  $\delta_i = 0$  for  $i = 1, \dots, b$ .

**for**  $s = 1, \dots, t$  **do**

**for all**  $(x_i, y_i) \in B$ : **do**

        Update  $\delta_i \leftarrow (1 - \eta_{\text{perturb}}\lambda)\delta_i + \eta_{\text{perturb}}\nabla_{\delta} \ell(\text{FORWARDPERTURB}(x_i, \delta_i, \Theta), y_i)$

Set update  $g = \nabla_{\Theta} [\frac{1}{b} \sum_i \ell(\text{FORWARDPERTURB}(x_i, \delta_i, \Theta), y_i)]$ .

Update  $\Theta \leftarrow \Theta - \eta(g + \nabla_{\Theta} R(\Theta))$ .

▷  $R$  is a regularizer, i.e. weight decay.

**function** FORWARDPERTURB( $x, \delta, \Theta$ )

▷ The net has layers  $f_1(\cdot; \Theta), \dots, f_r(\cdot; \Theta)$ ,  
with intended perturbations  $\delta^{(1)}, \dots, \delta^{(r)}$ .

Initialize  $h \leftarrow x$ .

**for**  $j = 1, \dots, r$  **do**

    Update  $h \leftarrow f_j(h; \Theta)$ .

    Update  $h \leftarrow h + \|h\|\delta^{(j)}$ .

**return**  $h$

---

Table 1: Validation error on CIFAR for standard training vs. AMO (Algorithm 1).

Dataset	Arch.	Setting	Standard SGD	AMO
CIFAR-10	WRN16-10	Baseline	4.15%	<b>3.42%</b>
		No data augmentation	9.59%	<b>6.74%</b>
		20% random labels	9.43%	<b>6.72%</b>
	WRN28-10	Baseline	3.82%	<b>3.00%</b>
		No data augmentation	8.28%	<b>6.47%</b>
		20% random labels	8.17%	<b>6.01%</b>
CIFAR-100	WRN16-10	Baseline	20.12%	<b>19.14%</b>
		No data augmentation	31.94%	<b>26.09%</b>
	WRN28-10	Baseline	18.85%	<b>17.78%</b>
		No data augmentation	30.04%	<b>24.67%</b>