### Generalization in Deep Neural Networks

Peter Bartlett

UC Berkeley Computer Science and Statistics

October 5, 2017

### • What determines the statistical complexity of a deep network?

- VC theory: Number of parameters
- Margins analysis: Size of parameters
- Understanding generalization failures

### Neural Networks for Classification

Neural network computes  $f : \mathbb{R}^d \to \mathbb{R}$ .

Directed acyclic graph with one output node, nodes compute

$$(z_1,\ldots,z_m)\mapsto\sigma\left(\sum_{i=1}^m w_iz_i+w_0\right),$$

where  $\sigma:\mathbb{R}\to\mathbb{R}$  is a nonlinear function, such as

$$\sigma(\alpha) = \frac{1}{1 + e^{-\alpha}},$$
 or  $\sigma(\alpha) = \max\{0, \alpha\}.$ 

Parameters (w's) adjusted by gradient descent to minimize an objective function on training examples, such as

$$\sum_{i=1}^n (f(x_i)-y_i)^2.$$

- Assume network maps to  $\{-1,1\}$ . (Threshold its output)
- Data generated by a probability distribution P on  $X \times \{-1, 1\}$ .
- Want to choose a function f such that with high probability  $P(f(x) \neq y)$  is small (near optimal).

# VC Theory

### Theorem (Vapnik and Chervonenkis)

Suppose  $F \subseteq \{-1,1\}^X$ . For every prob distribution P on  $X \times \{-1,1\}$ , with probability  $1 - \delta$  over n iid examples  $(x_1, y_1), \ldots, (x_n, y_n)$ , every f in F satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left( \operatorname{VCdim}(F) + \log(1/\delta) \right) \right)^{1/2}.$$

- For uniform bounds (that is, for all distributions and all *f* ∈ *F*, proportions are close to probabilities), this inequality is tight within a constant factor.
- For neural networks, VC-dimension:
  - increases with number of parameters
  - depends on nonlinearity and depth

#### Theorem

Consider the class F of  $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Piecewise constant (linear threshold units):

Piecewise linear (ReLUs):

Piecewise polynomial:

$$\operatorname{VCdim}(F) = \tilde{O}(p).$$

(Baum and Haussler, 1989)

 $\operatorname{VCdim}(F) = \tilde{O}(pL).$ 

(B., Harvey, Liaw, Mehrabian, 2017)

 $\operatorname{VCdim}(F) = \tilde{O}(pL^2).$ 

(B., Maiorov, Meir, 1998)

$$\operatorname{VCdim}(F) = \tilde{O}(p^2k^2).$$

(Karpinsky and MacIntyre, 1994)

### Sigmoid:

## Generalization in Neural Networks: Number of Parameters



### • What determines the statistical complexity of a deep network?

- VC theory: Number of parameters
- Margins analysis: Size of parameters
- Understanding generalization failures

### Theorem (B., 1996)

1. With high probability over *n* training examples  $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$ , every  $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$  has  $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$ 

2. If functions in F are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is,  $||w||_1 \leq B$ , then

$$\operatorname{fat}_{F}(\gamma) = \tilde{O}((B/\gamma)^{2}).$$

- The bound depends on the margin loss plus an error term.
- Minimizing quadratic loss or cross-entropy loss leads to large margins.
- fat<sub>F</sub>(γ) is a scale-sensitive version of VC-dimension. Unlike the VC-dimension, it need not grow with the number of parameters.



simons.berkeley.edu

- Qualitative behavior explained by small weights theorem.
- How to measure the complexity of a ReLU network?

### • What determines the statistical complexity of a deep network?

- VC theory: Number of parameters
- Margins analysis: Size of parameters
- Understanding generalization failures

### CIFAR10



http://corochann.com/

### Stochastic Gradient Training Error on CIFAR10



(Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, 2017) 13/24

### Training margins on CIFAR10 with true and random labels



- How does this match the large margin explanation?
- Need to account for the scale of the neural network functions.
- What is the appropriate notion of the size of these functions?

Spectrally-normalized margin bounds for neural networks. B., Dylan J. Foster, Matus Telgarsky, 2017. arXiv:1706.08498



Matus Telgarsky UIUC



Dylan Foster Cornell

## Generalization in Deep Networks

### New results for generalization in deep ReLU networks

- Measuring the size of functions computed by a network of ReLUs. (c.f. sigmoid networks: the output y of a layer has  $||y||_{\infty} \le 1$ , so  $||w||_1 \le B$  keeps the scale under control.)
- Large multiclass versus binary classification.

### Definitions

• Consider operator norms: For a matrix  $A_i$ ,

$$||A_i||_* := \sup_{||x|| \le 1} ||A_ix||.$$

• Multiclass margin function for  $f : \mathcal{X} \to \mathbb{R}^m$ ,  $y \in \{1, \dots, m\}$ :

$$M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i.$$

#### Theorem

With high probability, every  $f_A$  with  $R_A \leq r$  satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

### Definitions

Network with *L* layers, parameters  $A_1, \ldots, A_L$ :

$$f_A(x) := \sigma(A_L \sigma_{L-1}(A_{L-1} \cdots \sigma_1(A_1 x) \cdots)).$$

Scale of  $f_A$ :  $R_A := \prod_{i=1}^{L} \|A_i\|_* \sqrt{\sum_{i=1}^{L} \frac{\|A_i\|_F}{\|A_i\|_*}}$ .

(Assume  $\sigma_i$  is 1-Lipschitz, inputs normalized.)

### Stochastic Gradient Training Error on CIFAR10



(Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, 2017) 18/24

### Training margins on CIFAR10 with true and random labels



• How does this match the large margin explanation?

If we rescale the margins by  $R_A$  (the scale parameter):



If we rescale the margins by  $R_A$  (the scale parameter):

### Rescaled margins on MNIST



#### Theorem

With high probability, every  $f_A$  with  $R_A \leq r$  satisfies

$$\Pr(M(f_{\mathcal{A}}(X), Y) \leq 0) \leq \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}[M(f_{\mathcal{A}}(X_i), Y_i) \leq \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Network with *L* layers, parameters  $A_1, \ldots, A_L$ :

$$f_A(x) := \sigma(A_L \sigma_{L-1}(A_{L-1} \cdots \sigma_1(A_1 x) \cdots)).$$

Scale of  $f_A$ :  $R_A := \prod_{i=1}^{L} \|A_i\|_* \sqrt{\sum_{i=1}^{L} \frac{\|A_i\|_F}{\|A_i\|_*}}$ .



- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Lower bounds?
- Regularization: explicit control of operator norms?
- Role of depth?
- Interplay with optimization?
- Residual networks: Close to identity.