

Some Statistical Properties of Deep Networks

Peter Bartlett

UC Berkeley

August 2, 2018

Deep compositions of nonlinear functions

$$h = h_m \circ h_{m-1} \circ \cdots \circ h_1$$

Deep compositions of nonlinear functions

$$h = h_m \circ h_{m-1} \circ \cdots \circ h_1$$

e.g., $h_i : x \mapsto \sigma(W_i x)$

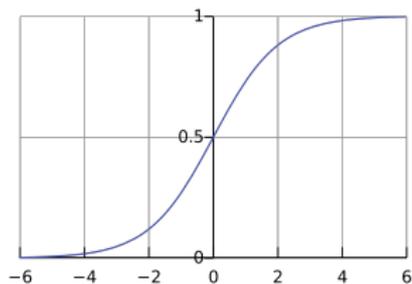
$$\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$$

Deep compositions of nonlinear functions

$$h = h_m \circ h_{m-1} \circ \cdots \circ h_1$$

e.g., $h_i : x \mapsto \sigma(W_i x)$

$$\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$$



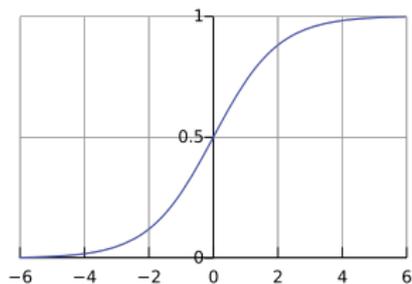
Deep compositions of nonlinear functions

$$h = h_m \circ h_{m-1} \circ \cdots \circ h_1$$

e.g., $h_i : x \mapsto \sigma(W_i x)$

$$\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$$

$$h_i : x \mapsto r(W_i x)$$
$$r(v)_i = \max\{0, v_i\}$$



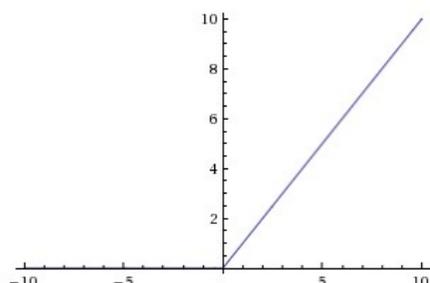
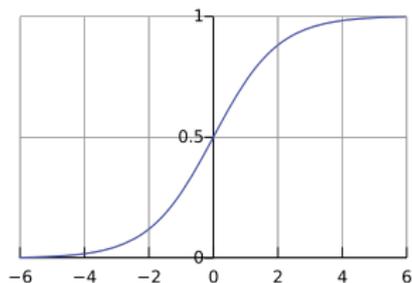
Deep compositions of nonlinear functions

$$h = h_m \circ h_{m-1} \circ \cdots \circ h_1$$

e.g., $h_i : x \mapsto \sigma(W_i x)$

$$\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$$

$h_i : x \mapsto r(W_i x)$
$$r(v)_i = \max\{0, v_i\}$$



Representation learning

Rich non-parametric family

Representation learning

Depth provides an effective way of representing useful features.

Rich non-parametric family

Representation learning

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representations.

Nonlinear parameterizations provide better rates of approximation.

(Birman & Solomjak, 1967), (DeVore et al, 1991)

Representation learning

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representations.

Nonlinear parameterizations provide better rates of approximation.

(Birman & Solomjak, 1967), (DeVore et al, 1991)

Some functions require much more complexity for a shallow representation.

(Telgarsky, 2015), (Eldan & Shamir, 2015)

Representation learning

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representations.

Nonlinear parameterizations provide better rates of approximation.

(Birman & Solomjak, 1967), (DeVore et al, 1991)

Some functions require much more complexity for a shallow representation.

(Telgarsky, 2015), (Eldan & Shamir, 2015)

- Statistical complexity?

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

- Assume network maps to $\{-1, 1\}$.
(Threshold its output)

- Assume network maps to $\{-1, 1\}$.
(Threshold its output)
- Data generated by a probability distribution P on $\mathcal{X} \times \{-1, 1\}$.

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

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$.

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

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

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$.

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

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

- For uniform bounds (that is, for all distributions and all $f \in \mathcal{F}$, proportions are close to probabilities), this inequality is tight—within a constant factor.

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$.

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

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

- For uniform bounds (that is, for all distributions and all $f \in \mathcal{F}$, proportions are close to probabilities), this inequality is tight—within a constant factor.
- For neural networks, VC-dimension:

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$.

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

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

- For uniform bounds (that is, for all distributions and all $f \in \mathcal{F}$, proportions are close to probabilities), this inequality is tight—within a constant factor.
- For neural networks, VC-dimension:
 - increases with number of parameters

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$.

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

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

- For uniform bounds (that is, for all distributions and all $f \in \mathcal{F}$, 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 \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Theorem

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

- 1 Piecewise constant (linear threshold units):

$$\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(p).$$

(Baum and Haussler, 1989)

Theorem

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

- 1 Piecewise constant (linear threshold units): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(p)$.
(Baum and Haussler, 1989)
- 2 Piecewise linear (ReLU): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(pL)$.
(B., Harvey, Liaw, Mehrabian, 2017)

Theorem

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

- 1 Piecewise constant (linear threshold units): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(p)$.
(Baum and Haussler, 1989)
- 2 Piecewise linear (ReLU): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(pL)$.
(B., Harvey, Liaw, Mehrabian, 2017)
- 3 Piecewise polynomial: $\text{VCdim}(\mathcal{F}) = \tilde{O}(pL^2)$.
(B., Maierov, Meir, 1998)

Theorem

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

- 1 Piecewise constant (linear threshold units): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(p)$.
(Baum and Haussler, 1989)
- 2 Piecewise linear (ReLU): $\text{VCdim}(\mathcal{F}) = \tilde{\Theta}(pL)$.
(B., Harvey, Liaw, Mehrabian, 2017)
- 3 Piecewise polynomial: $\text{VCdim}(\mathcal{F}) = \tilde{\mathcal{O}}(pL^2)$.
(B., Maierov, Meir, 1998)
- 4 Sigmoid: $\text{VCdim}(\mathcal{F}) = \tilde{\mathcal{O}}(p^2k^2)$.
(Karpinsky and MacIntyre, 1994)

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

Generalization in Deep Networks

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



Dylan Foster
Cornell



Matus Telgarsky
UIUC

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.
- For a pattern-label pair $(x, y) \in \mathcal{X} \times \{1, \dots, m\}$, define the margin $M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i$.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.
- For a pattern-label pair $(x, y) \in \mathcal{X} \times \{1, \dots, m\}$, define the margin $M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i$.
- If $M(f(x), y) > 0$ then f classifies x correctly.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.
- For a pattern-label pair $(x, y) \in \mathcal{X} \times \{1, \dots, m\}$, define the margin $M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i$.
- If $M(f(x), y) > 0$ then f classifies x correctly.
- We can view a larger margin as a more confident correct classification.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.
- For a pattern-label pair $(x, y) \in \mathcal{X} \times \{1, \dots, m\}$, define the margin $M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i$.
- If $M(f(x), y) > 0$ then f classifies x correctly.
- We can view a larger margin as a more confident correct classification.
- Minimizing a continuous loss, such as

$$\sum_{i=1}^n \|f(X_i) - Y_i\|^2,$$

encourages large margins.

Large-Margin Classifiers

- Consider a vector-valued function $f : \mathcal{X} \rightarrow \mathbb{R}^m$ used for classification, $y \in \{1, \dots, m\}$.
- The prediction on $x \in \mathcal{X}$ is $\arg \max_y f(x)_y$.
- For a pattern-label pair $(x, y) \in \mathcal{X} \times \{1, \dots, m\}$, define the margin $M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i$.
- If $M(f(x), y) > 0$ then f classifies x correctly.
- We can view a larger margin as a more confident correct classification.
- Minimizing a continuous loss, such as

$$\sum_{i=1}^n \|f(X_i) - Y_i\|^2,$$

encourages large margins.

- For large-margin classifiers, we should expect the fine-grained details of f to be less important.

New results for generalization in deep ReLU networks

- Measure the size of functions computed by a network of ReLUs via operator norms.

New results for generalization in deep ReLU networks

- Measure the size of functions computed by a network of ReLUs via operator norms.
- Large multiclass versus binary classification.

Generalization in Deep Networks

New results for generalization in deep ReLU networks

- Measure the size of functions computed by a network of ReLUs via operator norms.
- Large multiclass versus binary classification.

Definitions

- Consider operator norms: For a matrix A_i ,

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

Generalization in Deep Networks

New results for generalization in deep ReLU networks

- Measure the size of functions computed by a network of ReLUs via operator norms.
- Large multiclass versus binary classification.

Definitions

- Consider operator norms: For a matrix A_i ,
$$\|A_i\|_* := \sup_{\|x\| \leq 1} \|A_i x\|.$$
- Recall: Multiclass margin function for $f : \mathcal{X} \rightarrow \mathbb{R}^m$, $y \in \{1, \dots, m\}$, is

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

Generalization in Deep Networks

Theorem

With high probability, every f_A

Generalization in Deep Networks

Theorem

With high probability, every f_A

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Generalization in Deep Networks

Theorem

With high probability, every f_A satisfies

$$\Pr(M(f_A(X), Y) \leq 0) \leq$$

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Generalization in Deep Networks

Theorem

With high probability, every f_A satisfies

$$\Pr(M(f_A(X), Y) \leq 0) \leq \frac{1}{n} \sum_{i=1}^n 1[M(f_A(X_i), Y_i) \leq \gamma]$$

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Generalization in Deep Networks

Theorem

With high probability, every f_A satisfies

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

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Generalization in Deep Networks

Theorem

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

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

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Scale of f_A : $R_A := \prod_{i=1}^L \|A_i\|_*$

Generalization in Deep Networks

Theorem

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

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

Definitions

Network with L layers, parameters A_1, \dots, A_L :

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

Scale of f_A : $R_A := \prod_{i=1}^L \|A_i\|_* \left(\sum_{i=1}^L \frac{\|A_i\|_{2,1}^{2/3}}{\|A_i\|_*^{2/3}} \right)^{3/2}$.

(Assume σ_i is 1-Lipschitz, inputs normalized.)

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

Understanding Generalization Failures

CIFAR10

6: frog



9: truck



9: truck



4: deer



1: automobile



1: automobile



2: bird



7: horse



8: ship



3: cat



4: deer



7: horse



7: horse



2: bird



9: truck



9: truck



9: truck



3: cat



2: bird



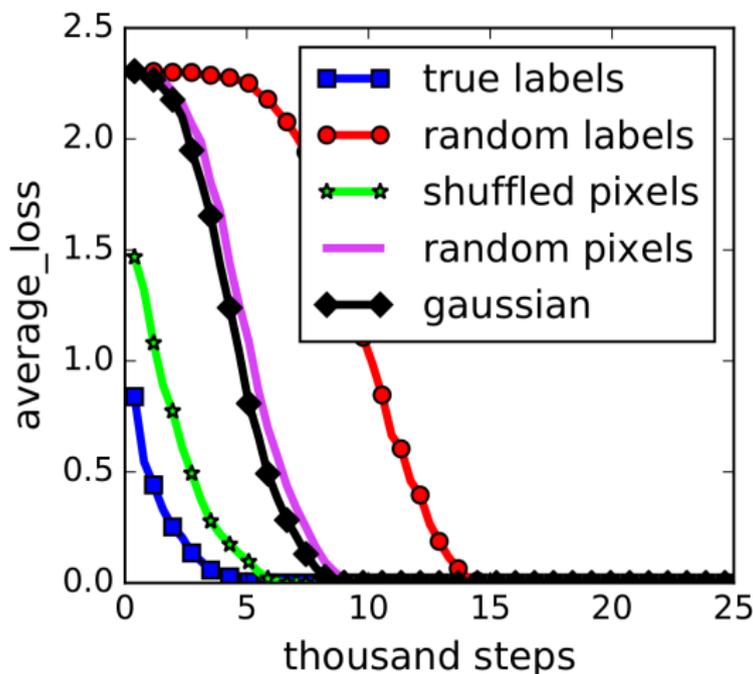
6: frog



<http://corochann.com/>

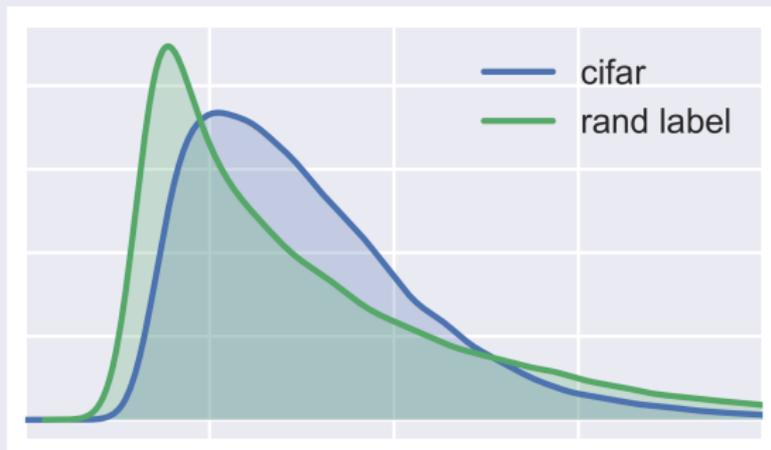
Understanding Generalization Failures

Stochastic Gradient Training Error on CIFAR10



Understanding Generalization Failures

Training margins on CIFAR10 with true and random labels

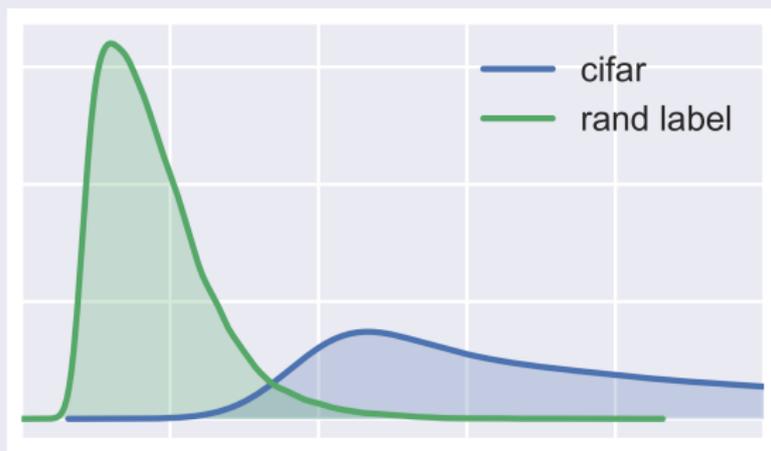


- How does this match the large margin explanation?

Understanding Generalization Failures

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

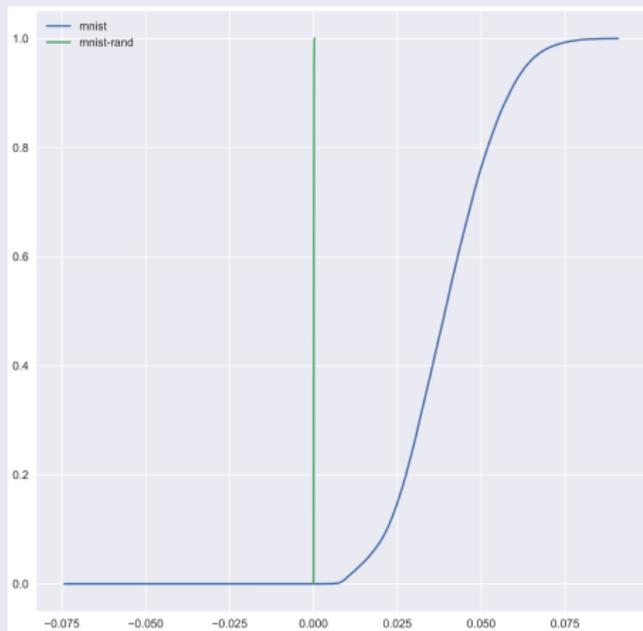
Rescaled margins on CIFAR10



Understanding Generalization Failures

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

Rescaled cumulative margins on MNIST



Generalization in Deep Networks

Theorem

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

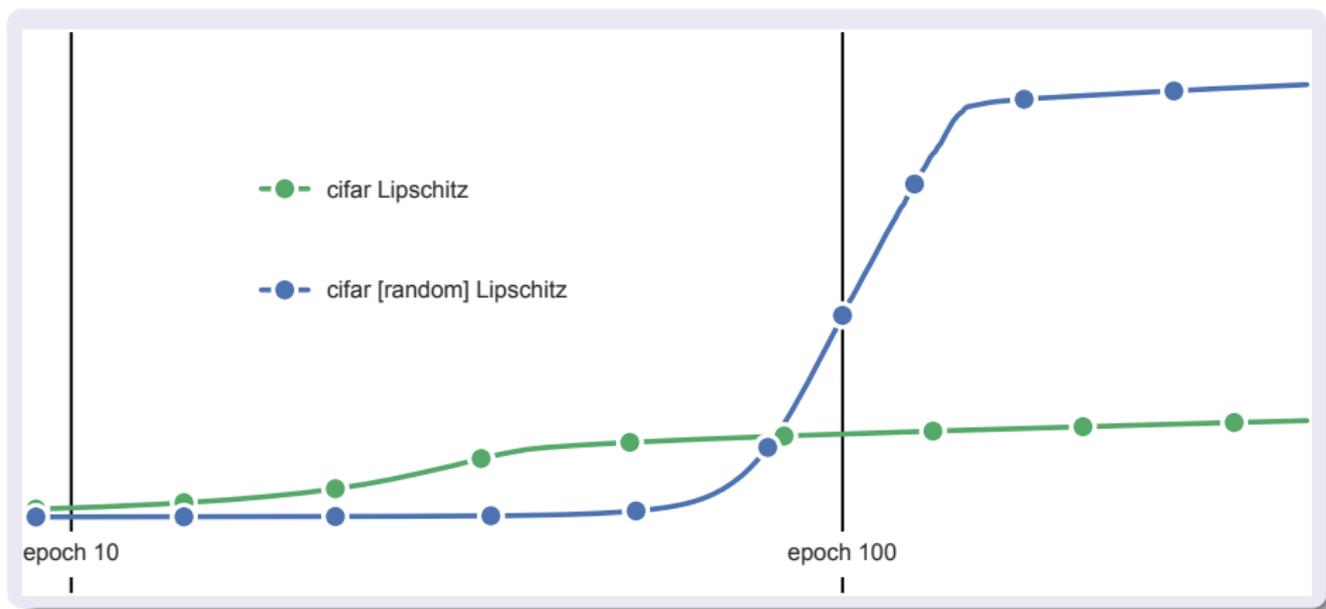
$$\Pr(M(f_A(X), Y) \leq 0) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \leq \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Network with L layers, parameters A_1, \dots, A_L :

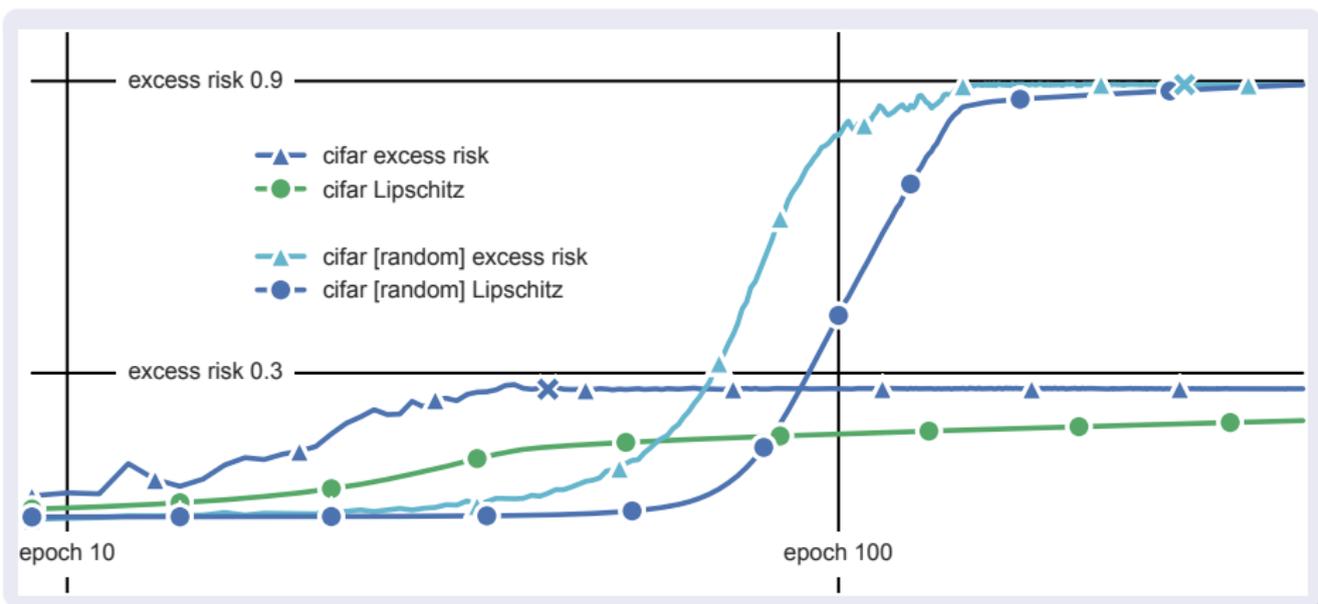
$$f_A(x) := \sigma(A_L \sigma_{L-1}(A_{L-1} \cdots \sigma_1(A_1 x) \cdots)).$$

Scale of f_A : $R_A := \prod_{i=1}^L \|A_i\|_* \left(\sum_{i=1}^L \frac{\|A_i\|_{2,1}^{2/3}}{\|A_i\|_*^{2/3}} \right)^{3/2}$.

Understanding Generalization Failures



Understanding Generalization Failures



- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.

Generalization in Neural Networks

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.

Generalization in Neural Networks

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.
- Regularization and optimization: explicit control of operator norms?