

# The Role of Over-Parametrization in Generalization of Neural Networks

Behnam Neyshabur  
NYU

Zhiyuan Li  
Princeton

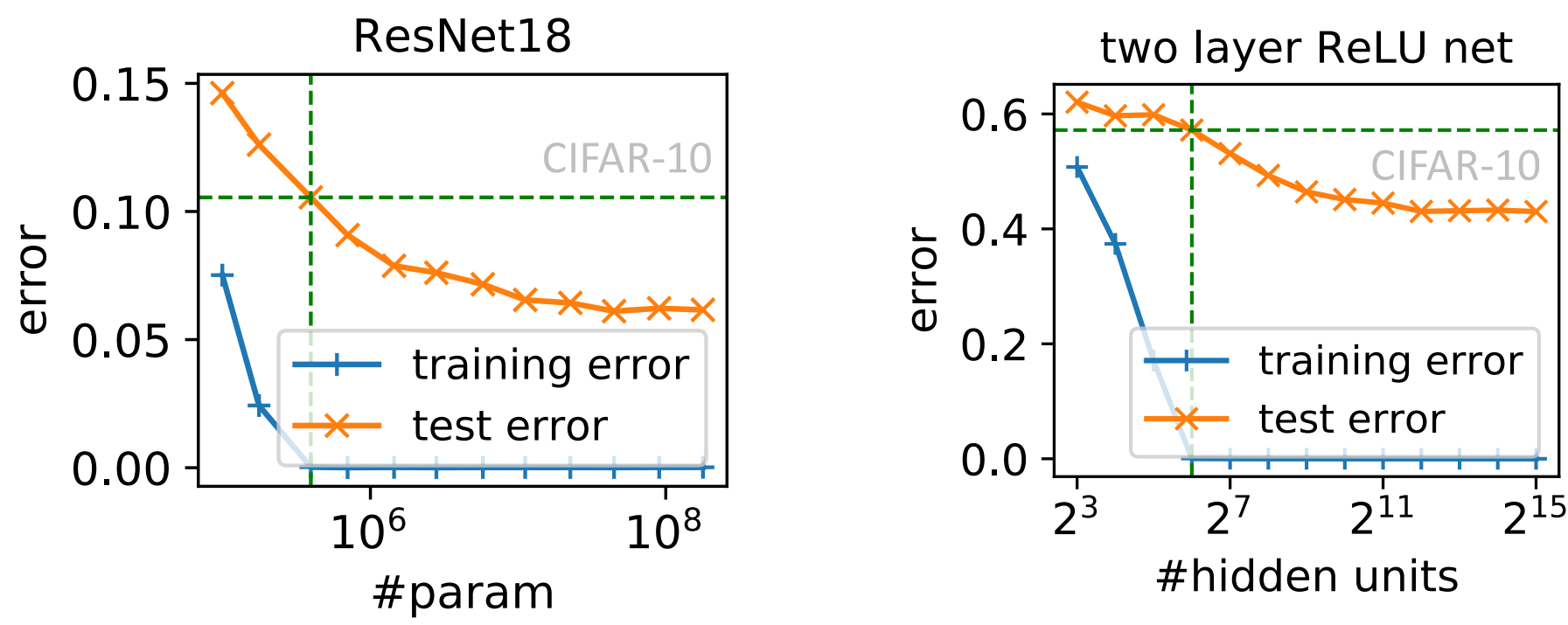
Srinadh Bhojanapalli  
Google

Yann LeCun  
NYU

Nathan Srebro  
TTI-Chicago

Empirical observation:

Over-parametrization helps generalization



Current complexity measures  $\uparrow$  with over-parametrization  $\text{☹}$

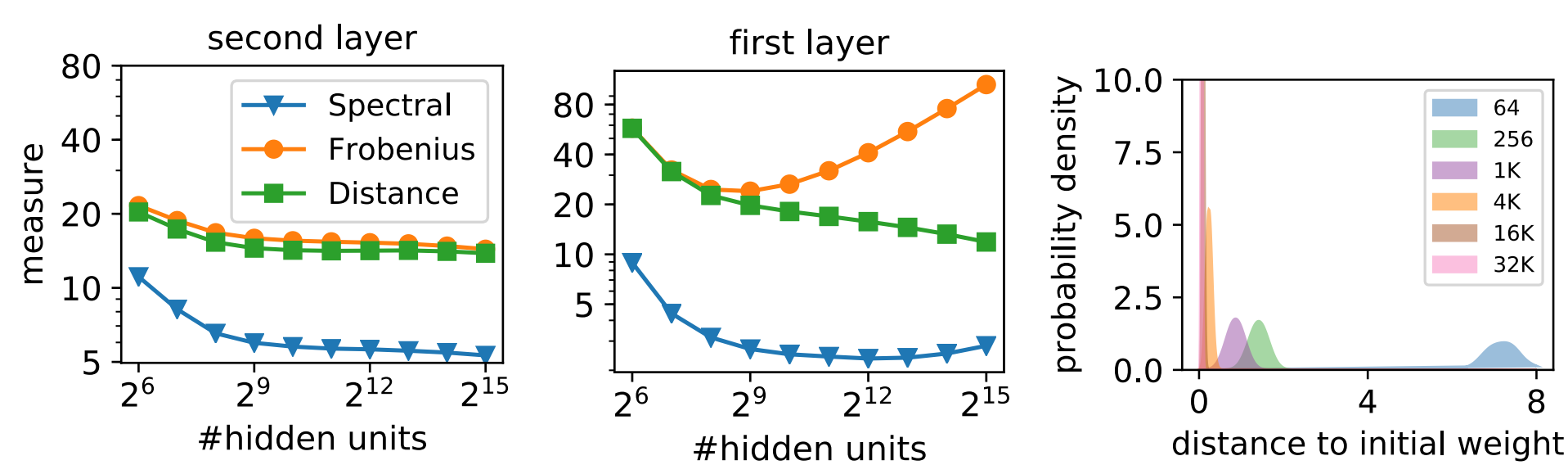
A generalization bound that  $\downarrow$  with over-parametrization?  $\text{☺}$

Why is this important and useful?

- One of the key **mysteries** behind success of deep learning
- Understanding the real cause of generalization can be used to **design better models and optimization algorithms**

## An Empirical Investigation

Properties of learned two layer ReLU networks:

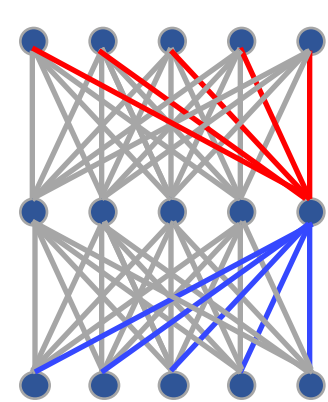
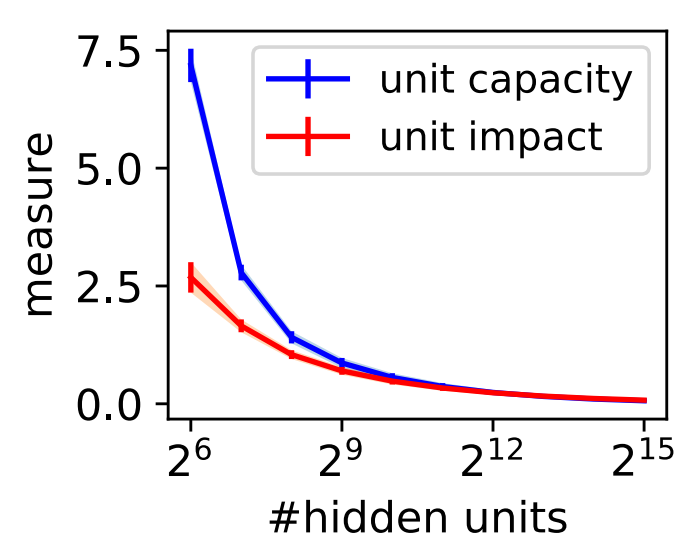


Distance: Euclidean distance from initialization.

$$\mathcal{F}_{\mathcal{W}} = \{f(\mathbf{x}) = \mathbf{V}[\mathbf{U}\mathbf{x}]_+ \mid (\mathbf{V}, \mathbf{U}) \in \mathcal{W}\}$$

$$\mathcal{W} = \{(\mathbf{V}, \mathbf{U}) \mid \mathbf{V} \in \mathbb{R}^{c \times h}, \mathbf{U} \in \mathbb{R}^{h \times d}, \|\mathbf{v}_i\| \leq \alpha_i, \|\mathbf{u}_i - \mathbf{u}_i^0\|_2 \leq \beta_i\}$$

unit impact  $\uparrow$  unit capacity  $\uparrow$



## Generalization Bound

Rademacher complexity bound:

$$\mathcal{R}_{\mathcal{S}}(\ell_{\gamma} \circ \mathcal{F}_{\mathcal{W}}) \leq \frac{2\sqrt{2c} + 2}{\gamma m} \sum_{j=1}^h \alpha_j \left( \beta_j \|\mathbf{X}\|_F + \|\mathbf{u}_j^0 \mathbf{X}\|_2 \right)$$

Generalization bound:

For any  $m$ , w.p.  $> 1 - \delta$  over the training set:

$$L_0(f) \leq \hat{L}_{\gamma}(f) + \tilde{O} \left( \frac{\sqrt{c} \|\mathbf{V}\|_F (\|\mathbf{U} - \mathbf{U}^0\|_F + \|\mathbf{U}^0\|_2) \sqrt{\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i\|_2^2}}{\gamma \sqrt{m}} + \sqrt{\frac{h}{m}} \right)$$

$c$ : #classes  $h$ : #hidden units  $d$ : input dim  $m$ : #samples  $\gamma$ : margin

•  $\mathbf{U}^0$  is the initialization / any fixed matrix.

• Rademacher complexity:  $\mathcal{R}_{\mathcal{S}}(\mathcal{H}) = \frac{1}{m} \mathbb{E} \left[ \sup_{f \in \mathcal{H}} \sum_{i=1}^m \xi_i f(x_i) \right]$

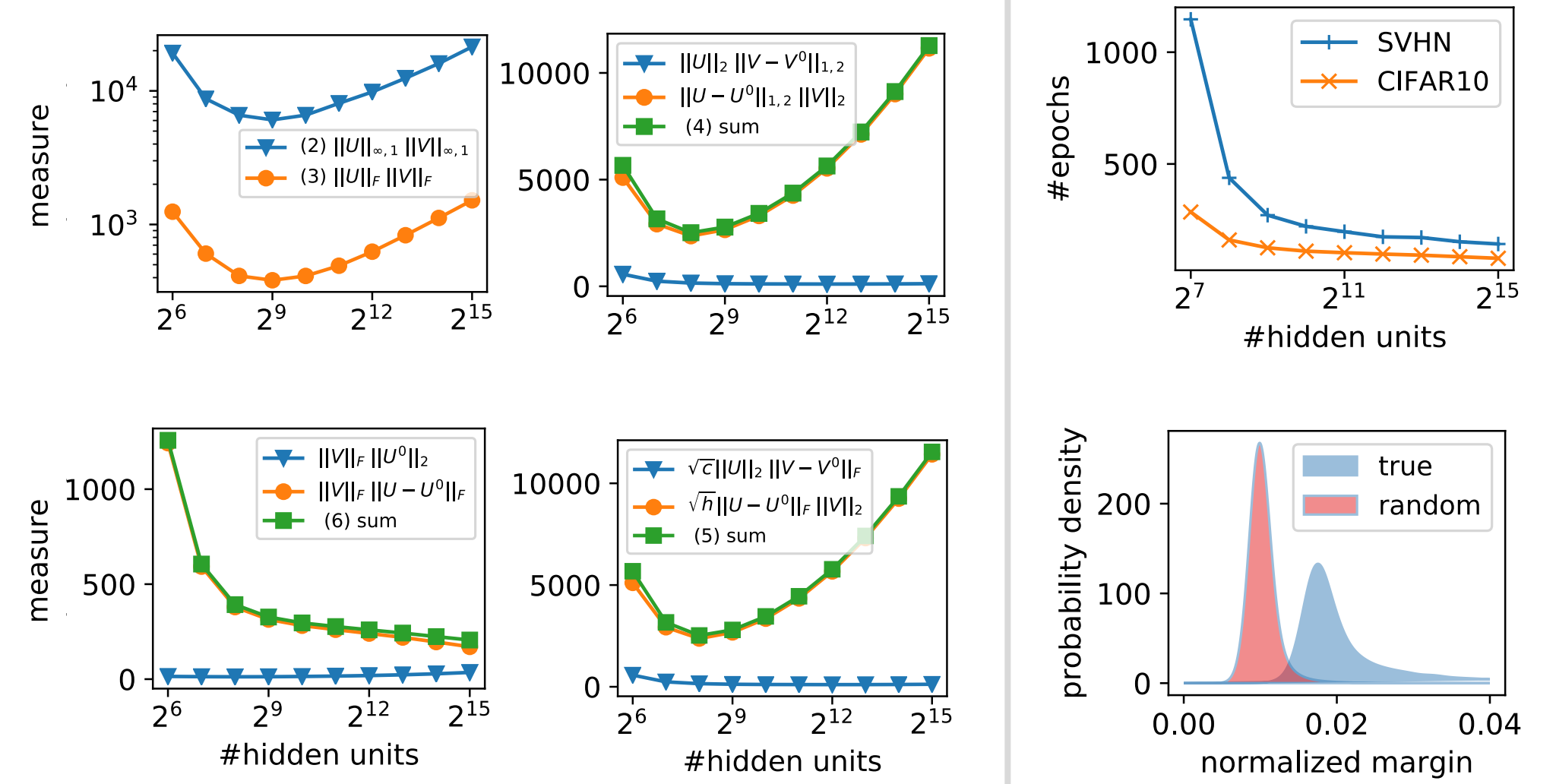
• Margin ramp loss:  $\ell_{\gamma}(f(\mathbf{x}), y) = \begin{cases} 0 & \mu(f(\mathbf{x}), y) > \gamma \\ \mu(f(\mathbf{x}), y) / \gamma & \mu(f(\mathbf{x}), y) \in [0, \gamma] \\ 1 & \mu(f(\mathbf{x}), y) < 0. \end{cases}$

## Evaluation

Comparing generalization measures:

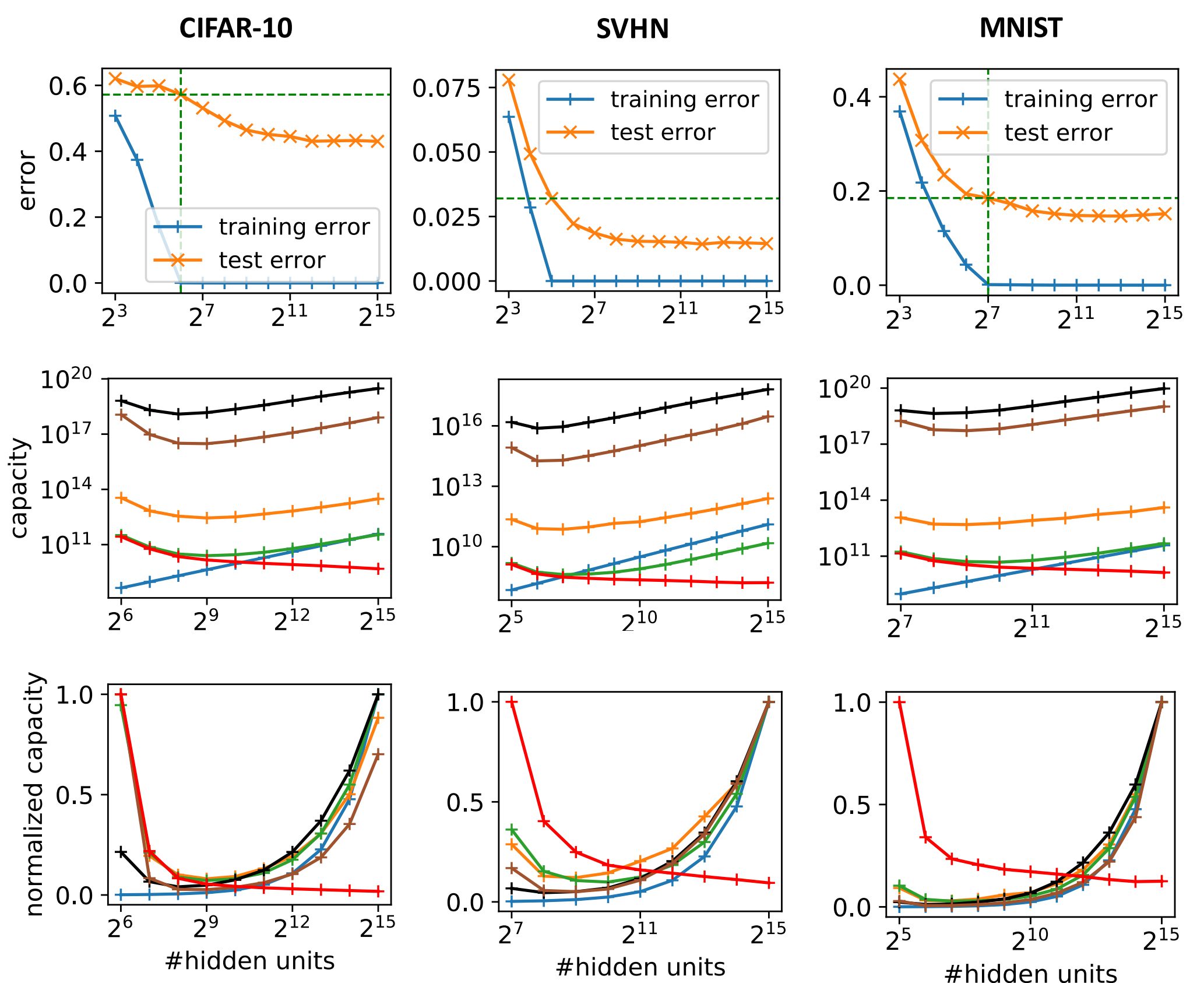
#	Reference	Measure
(1)	Harvey et al. 2017	$\tilde{O}(dh)$
(2)	Bartlett and Mendelson 2002	$\tilde{O}(\ \mathbf{U}\ _{\infty,1} \ \mathbf{V}\ _{\infty,1})$
(3)	Neyshabur et al. 2015	$\tilde{O}(\ \mathbf{U}\ _F \ \mathbf{V}\ _F)$
(4)	Bartlett et al. 2017	$\tilde{O}(\ \mathbf{U}\ _2 \ \mathbf{V} - \mathbf{V}_0\ _{1,2} + \ \mathbf{U} - \mathbf{U}_0\ _{1,2} \ \mathbf{V}\ _2)$
(5)	Neyshabur et al. 2018	$\tilde{O}(\ \mathbf{U}\ _2 \ \mathbf{V} - \mathbf{V}_0\ _F + \sqrt{h} \ \mathbf{U} - \mathbf{U}_0\ _F \ \mathbf{V}\ _2)$
(6)	ours	$\tilde{O}(\ \mathbf{U}_0\ _2 \ \mathbf{V}\ _F + \ \mathbf{U} - \mathbf{U}^0\ _F \ \mathbf{V}\ _F + \sqrt{h})$

Behavior of different norms with over-parametrization:



Comparing capacity bounds:

(1) VC-dim (2)  $\ell_{1,\infty}$  (3) Fro (4) spec- $\ell_{2,1}$  (5) spec-Fro (6) ours



Code available at <https://github.com/bneyshabur/over-parametrization>

## Matching Lower Bound

Consider  $\mathcal{W}' = \{(\mathbf{V}, \mathbf{U}) \mid \mathbf{V} \in \mathbb{R}^{1 \times h}, \mathbf{U} \in \mathbb{R}^{h \times d}, \|\mathbf{v}_j\| \leq \alpha_j, \|\mathbf{u}_j - \mathbf{u}_j^0\|_2 \leq \beta_j, \|\mathbf{U} - \mathbf{U}^0\|_2 \leq \max_{j \in [h]} \beta_j\}$

For any  $d = h \leq m$ ,  $\{\alpha_j, \beta_j\}_{j=1}^h \in \mathbb{R}^+$ ,  $\mathbf{U}_0 \in \mathbb{R}^{h \times d}$  there exists  $S = \{\mathbf{x}_i\}_{i=1}^m \subset \mathbb{R}^d$  s.t.

$$\mathcal{R}_{\mathcal{S}}(\mathcal{F}_{\mathcal{W}}) \geq \mathcal{R}_{\mathcal{S}}(\mathcal{F}_{\mathcal{W}'}) = \Omega \left( \frac{\sum_{j=1}^h \alpha_j \|\mathbf{u}_j^0 \mathbf{X}\|_2}{m} \right)$$

For any  $d = h \leq m$ ,  $\{\alpha_j, \beta_j\}_{j=1}^h \in \mathbb{R}^+$ ,  $\mathbf{U}_0 = 0$  there exists  $S = \{\mathbf{x}_i\}_{i=1}^m \subset \mathbb{R}^d$  s.t.

$$\mathcal{R}_{\mathcal{S}}(\mathcal{F}_{\mathcal{W}}) \geq \mathcal{R}_{\mathcal{S}}(\mathcal{F}_{\mathcal{W}'}) = \Omega \left( \frac{\sum_{j=1}^h \alpha_j \beta_j \|\mathbf{X}\|_F}{m} \right)$$

✓ Matches both terms in the upper bound.

✓ First lower bound that is higher than Lipschitz of the network.

✓ Improves by a factor of  $\sqrt{h}$  over Bartlett et al. 2017.

✓  $\sqrt{h}$  gap between linear vs non-linear networks.