

Generalization Theory for Deep Learning



Basic version: finite hypothesis class

Finite hypothesis class: with probability $1 - \delta$ over the choice of a training set of size n , for a bounded loss ℓ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)] \right| = O \left(\sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}} \right)$$

VC-Dimension

Motivation: Do we need to consider **every** classifier in \mathcal{F} ?

Intuitively, **pattern of classifications** on the training set should suffice. (Two predictors that predict identically on the training set should generalize similarly).

Let $\mathcal{F} = \{f : \mathbb{R}^d \rightarrow \{+1, -1\}\}$ be a class of binary classifiers.

The **growth function** $\Pi_{\mathcal{F}} : \mathbb{N} \rightarrow \mathbb{F}$ is defined as:

$$\Pi_{\mathcal{F}}(m) = \max_{(x_1, x_2, \dots, x_m)} \left| \left\{ (f(x_1), f(x_2), \dots, f(x_m)) \mid f \in \mathcal{F} \right\} \right|.$$

The **VC dimension** of \mathcal{F} is defined as:

$$\text{VCdim}(\mathcal{F}) = \max \{ m : \Pi_{\mathcal{F}}(m) = 2^m \}.$$

VC-dimension Generalization bound

Theorem (Vapnik-Chervonenkis): with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

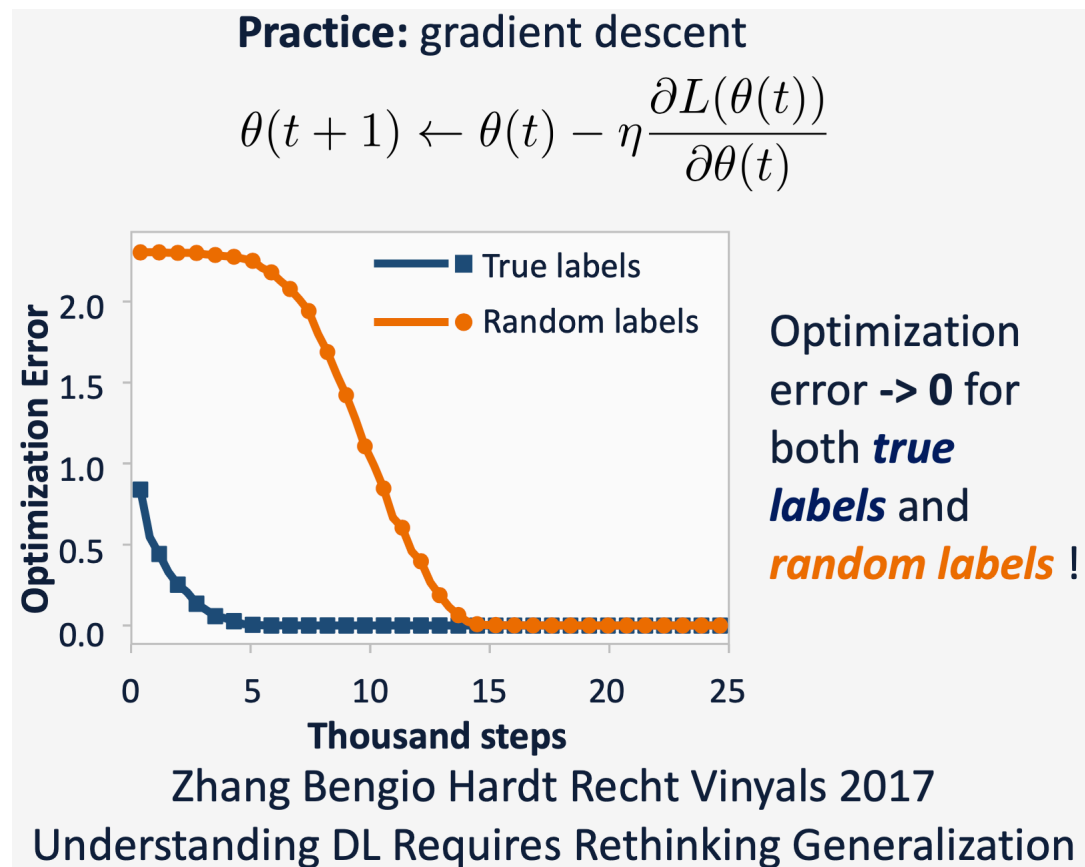
$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)] \right| = O \left(\sqrt{\frac{\text{VCdim}(\mathcal{F}) \log n + \log 1/\delta}{n}} \right)$$

Examples:

- Linear functions: VC-dim = $O(\text{dimension})$
- Neural network: VC-dimension of fully-connected net with width W and H layers is $\widetilde{\Theta}(WH)$ (Bartlett et al., '17).

Problems with VC-dimension bound

1. In over-parameterized regime, bound $\gg 1$.
2. Cannot explain the random noise phenomenon:
 - Neural networks that fit random labels and that fit true labels have the same VC-dimension.



PAC Bayesian Generalization Bounds

Setup: Let P be a prior over function in class \mathcal{F} , let Q be the posterior (after algorithm's training).

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)] \right| = O \left(\sqrt{\frac{KL(Q || P) + \log 1/\delta}{n}} \right)$$

$$P(W_{ij}) \sim \mathcal{N}(0, 1)$$

P does not depend on data

data-dependent bound

Rademacher Complexity

Intuition: how well can a classifier class fit random noise?

(Empirical) **Rademacher complexity:** For a training set $S = \{x_1, x_2, \dots, x_n\}$, and a class \mathcal{F} , denote:

$$\hat{R}_n(S) = \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(x_i).$$

$(\sigma_1, \dots, \sigma_n)$
random noise

where $\sigma_i \sim \text{Unif}\{+1, -1\}$ (Rademacher R.V.).

(Population) **Rademacher complexity:**

$$R_n = \mathbb{E}_S \left[\hat{R}_n(S) \right].$$

Rademacher Complexity Generalization Bound

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)] \right| = O \left(\frac{\hat{R}_n}{n} + \sqrt{\frac{\log 1/\delta}{n}} \right)$$

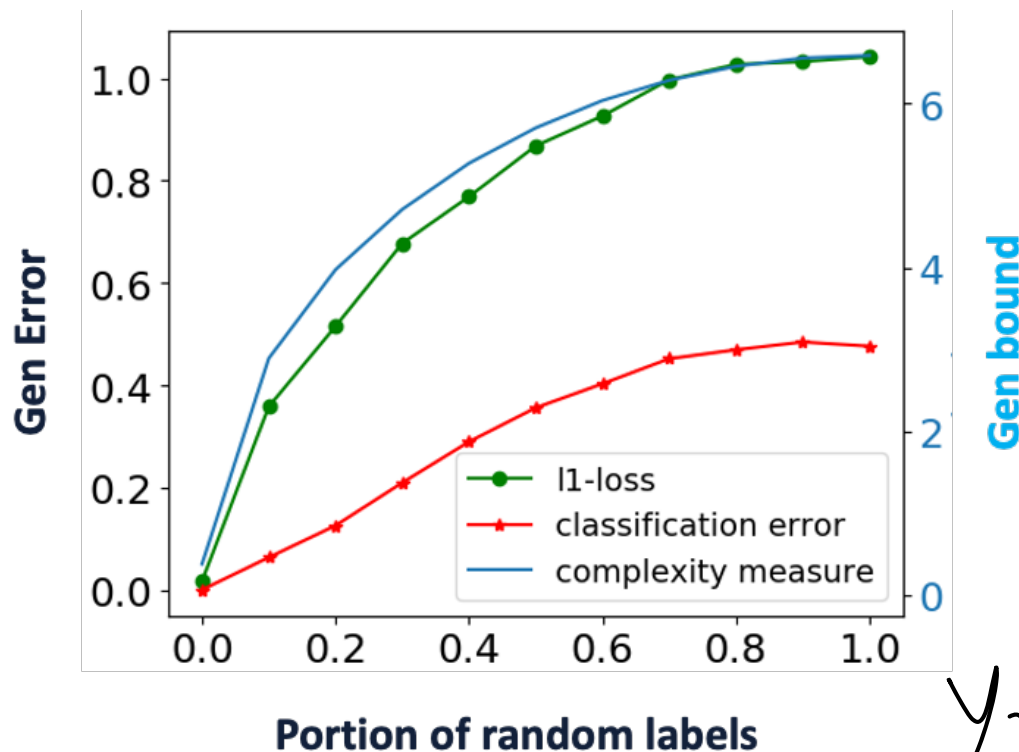
and

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)] \right| = O \left(\frac{R_n}{n} + \sqrt{\frac{\log 1/\delta}{n}} \right)$$

Kernel generalization bound

data-dependent
 H^*, y

Use Rademacher complexity theory, we can obtain a generalization bound $O(\sqrt{y^\top (H^*)^{-1} y / n})$ where $y \in \mathbb{R}^n$ are n labels, and $H^* \in \mathbb{R}^{n \times n}$ is the kernel (e.g., NTK) matrix.



$$y_i = (1-\alpha) \cdot y_i^{\text{true}} + \alpha \cdot \text{noise}$$

α

Norm-based Rademacher complexity bound

$$\underline{ReLU}: \rho = 1$$

Theorem: If the activation function σ is ρ -Lipschitz. Let

$$\mathcal{F} = \{x \mapsto W_{H+1} \sigma(W_H \sigma(\dots \sigma(W_1 x) \dots)), \|W_h^T\|_{1, \infty} \leq B \forall h \in [H]\}$$

\int : training data

then $R_n(\mathcal{S}) \leq \|X^T\|_{2, \infty} (2\rho B)^{H+1} \sqrt{2 \ln d}$ where

$X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ is the input data matrix.

$$\|X^T\|_{2, \infty} = \max_{i=1, \dots, d} \|x_i\|_2$$

$$\|W_h^T\|_{1, \infty} = \max_{\text{row}} \|W_h^{\text{row}}\|_1$$

Pf: induction layer by layer

Comments on generalization bounds

- When plugged in real values, the bounds are rarely non-trivial (i.e., smaller than 1)
- “*Fantastic Generalization Measures and Where to Find them*” by Jiang et al. '19 : large-scale investigation of the correlation of extant generalization measures with true generalization.

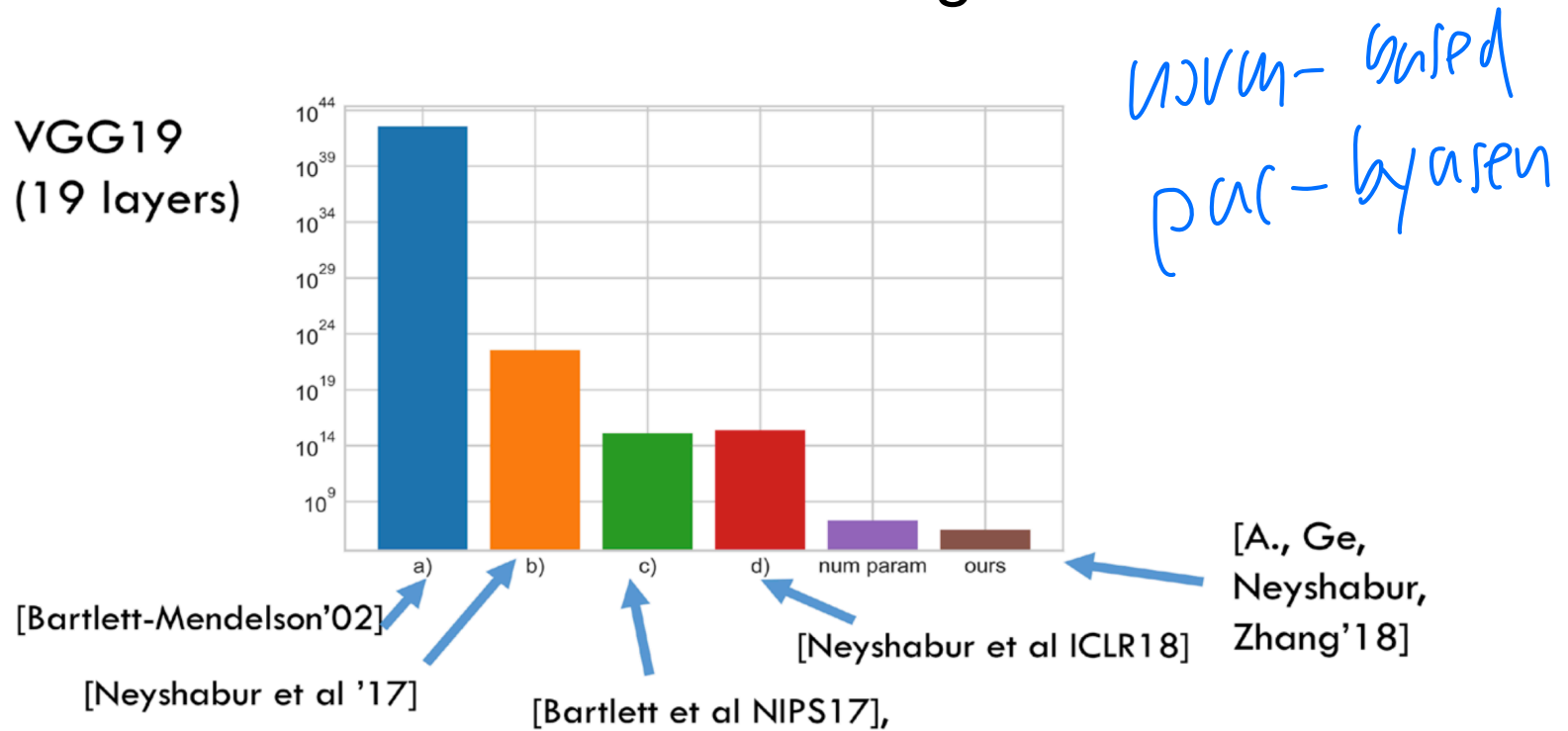


Image credits to Andrej Risteski

Comments on generalization bounds

- Uniform convergence may be unable to explain generalization of deep learning [Nagarajan and Kolter, '19]
 - Uniform convergence: a bound for all $f \in \mathcal{F}$
 - Exists example that 1) can generalize, 2) uniform convergence fails.

non-uniform: linear regression
bias² + Variance

- Rates:
 - Most bounds: $1/\sqrt{n}$.
 - Local Rademacher complexity: $1/n$.

Separation between NN and kernel

- For approximation and optimization, neural network has no advantage over kernel. Why NN gives better performance: **generalization.**
- [Allen-Zhu and Li '20] Construct a class of functions \mathcal{F} such that $y = f(x)$ for some $f \in \mathcal{F}$:
 - no kernel is sample-efficient; *need exponential # of data*
 - Exists a neural network that is sample-efficient. *polynomial samples*

Separation between NN and kernel

Defn Kernel method is a linear method with an embedding. $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$, Hilbert space

\Rightarrow it turns an element $f \in \mathcal{H}$ into a prediction function $y = \langle f, \phi(x) \rangle$

$\hat{y} = f(x)$

The method uses samples, $\{x_i\}_{i=1}^n$, $x_i \in \mathbb{R}^d$
observe $\{y_i\}_{i=1}^n$

$f \in \text{span} \left(\phi(x_i) \right)_{i=1}^n$, $i \in [n]$

e.g., $\underset{f}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \phi(x_i), f \rangle)^2 + \lambda \|f\|^2$

Separation between NN and kernel

Thm: \exists a class of functions $\mathcal{C} \subseteq \{c: \mathbb{R}^d \rightarrow \mathbb{R}\}$
and a distribution over \mathbb{R}^d s.t.

i) \forall kernel method, $\forall c \in \mathcal{C}$,
given $y_i = c(x_i)$
if $\mathbb{E}_{x \sim \mu} [(c(x) - \langle f, \phi(x) \rangle)^2] \leq \frac{1}{n}$
then $n \geq 2$ *exponential*

ii) \exists simple procedure, that can output
the true c as long as $n \geq d$

\star the procedure can be simulated by
a neural network f_{NN}

Separation between NN and kernel

Pf: μ : uniform distribution over $\{-1, 1\}^d$
 2^d element

$$\mathcal{C} = \left\{ C_S = \prod_{s \in S} X_s, \quad S \subset \{1, \dots, d\} \right\}$$

Pf of Part ii) choose a basis $\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ \vdots \\ -1 \end{pmatrix}, \dots, \begin{pmatrix} 1 \\ \vdots \\ -1 \end{pmatrix}$
 e_1, e_2, \dots, e_d

$$\Rightarrow y_i = C_S(p_i), \quad \text{if } i \in S, \quad y_i = -1$$

$$i \notin S, \quad y_i = 1$$

\Rightarrow Know whether i is in S or not

\Rightarrow identify $S \Rightarrow$ learn $C_S \square$

Separation between NN and kernel

Part i) \mathcal{C} is a basis for $\{f: \{+1,1\}^d \rightarrow \mathbb{R}\}$
with distribution μ .

$$\mathbb{E}_{x \sim \mu} [C_f(x) \cdot C_{f'}(x)] = \begin{cases} 0 & \text{if } f \neq f' \\ 1 & \text{if } f = f' \end{cases}$$

Goal: $\mathbb{E}_{x \sim \mu} [(C_{f^*}(x) - \langle f, \phi(x) \rangle)^2]$

1) since $f \in \text{span}(\phi(x_i))_{i=1}^n$
 $f = \sum_{i=1}^n a_i \phi(x_i)$, $f(x) = \sum_{i=1}^n a_i \langle \phi(x_i), \phi(x) \rangle$

$$x \mapsto \langle \phi(x_i), \phi(x) \rangle$$

$$= \sum_{f \in \mathcal{C}[d]} \lambda_i f(C_f(x))$$

Separation between NN and kernel

$$\begin{aligned} & \mathbb{E}_{X \sim \mathcal{M}} \left[\left(C_{S^*}(X) - \langle f, \phi(X) \rangle \right)^2 \right] \\ &= \mathbb{E}_{X \sim \mathcal{M}} \left[\left(C_{S^*}(X) - \sum_{S \in \mathcal{C}(d)} \sum_{i=1}^n a_i \cdot \lambda_{i,S} C_S(X) \right)^2 \right] \\ &= \left(1 - \sum_{i=1}^n a_i \lambda_{i,S^*} \right)^2 + \sum_{S \neq S^*} \left(\sum_{i=1}^n a_i \lambda_{i,S} \right)^2 \\ & \text{by assumption error} \leq \frac{1}{9} \\ & \Rightarrow \left(1 - \sum_{i=1}^n a_i \lambda_{i,S^*} \right)^2 \leq \frac{1}{9} \\ & \& \sum_{S \neq S^*} \left(\sum_{i=1}^n a_i \lambda_{i,S} \right)^2 \leq \frac{1}{9} \end{aligned}$$

Separation between NN and kernel

Notations

$$\Lambda : 2^d \times \mathcal{U}$$

$$\Omega_{S^+, S^*}$$

$$\Lambda_{S^+, i} = \Lambda_{i, S^+}$$

$$= \sum_{j=1}^4 \Lambda_{i, S^+} \cdot \alpha_{i, S^+}^j$$

$$A = n \times 2^d$$

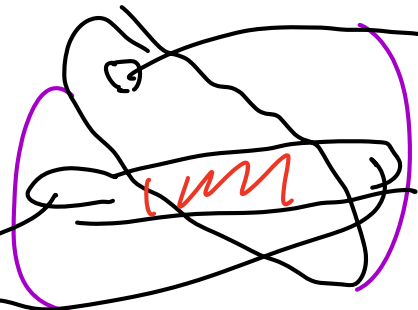
$$A_{i, S^*} = \alpha_{i, S^*}$$

$$\Omega = \Lambda A : 2^d \times 2^d \text{ of rank } \approx 4$$

$$\left(1 - \Omega_{S^+, S^*} \right)^2 \leq \frac{1}{9} \rightarrow \Omega_{S^+, S^*} \geq \frac{4}{9}$$

$$\sum_{S \neq S^*} \Omega_{S, S^*}^2 \leq \frac{1}{9}$$

$$\leq \frac{1}{9} \rightarrow \Omega : 2^d \times 2^d \Rightarrow \text{rank at least } \frac{3}{4} \cdot 2^d$$



Separation between NN and kernel

$$\Omega = \text{diag}(\Omega) + \Omega', \quad \Omega' : \text{off-diagonal}$$

$$\|\Omega'\|_F^2 = \frac{2^d}{4} = \sum \text{eigen}^2(\Omega)$$

$\Rightarrow \Omega'$ has at most $\frac{2^d}{4}$ eigenvalues

\Rightarrow consider subspace with eigenvalue $\geq \frac{2}{3}$ which has dimension at least $\frac{3}{4} \cdot 2^d$

$\forall x \in \text{subspace}$

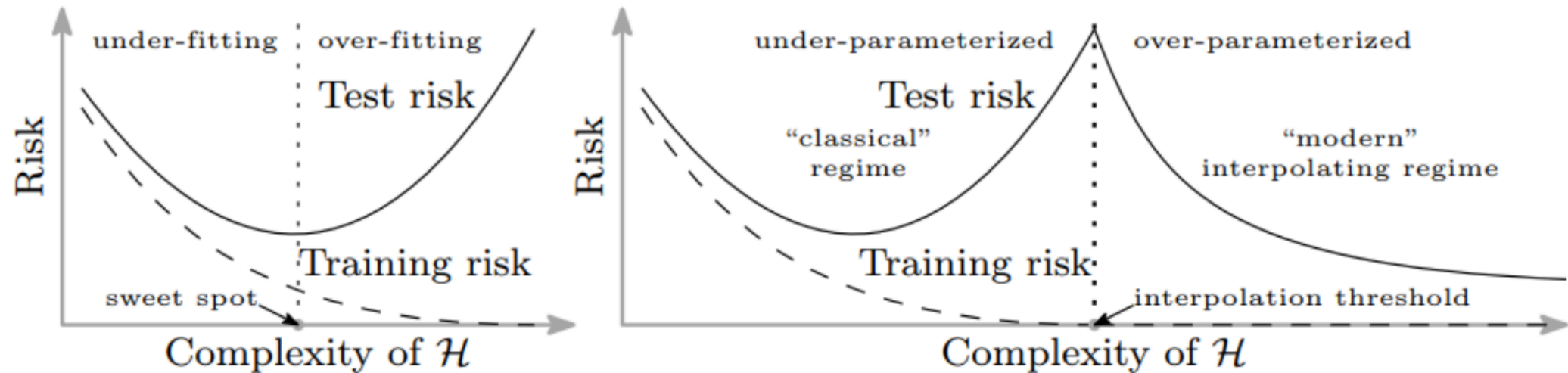
$$\|\Omega x\|_2 = \|\text{diag}(\Omega)x + \Omega'x\|_2$$

$\Rightarrow x \in \text{span}(\Omega)$

$\Rightarrow \text{span}(\Omega) \text{ dim } \geq \frac{3}{4} \cdot 2^d$

$$\geq \|\text{diag}(\Omega)x\|_2 - \|\Omega'x\|_2$$
$$\geq \frac{2}{3} \|x\| - \frac{2}{3} \|x\| = 0$$

Double descent



(a) U-shaped “bias-variance” risk curve

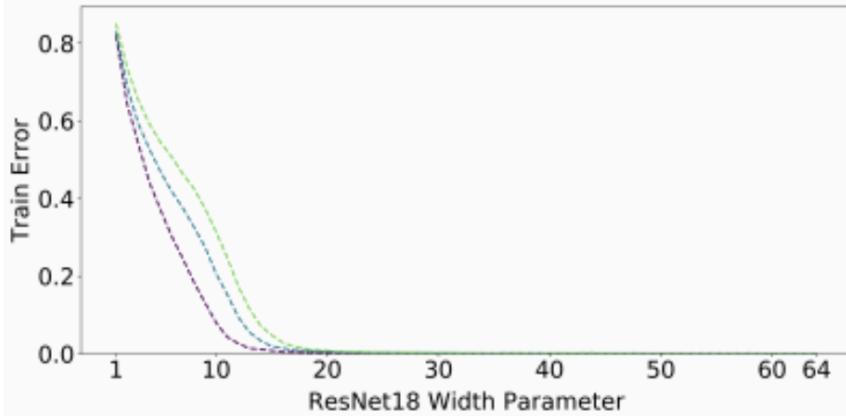
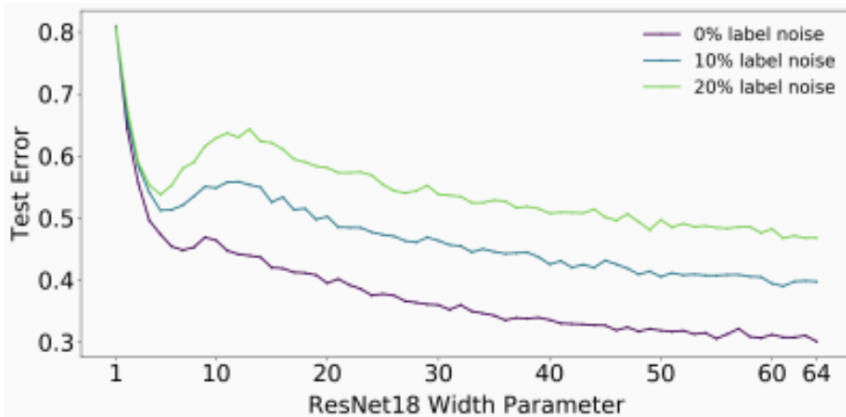
(b) “double descent” risk curve

Belkin, Hsu, Ma, Mandal '18

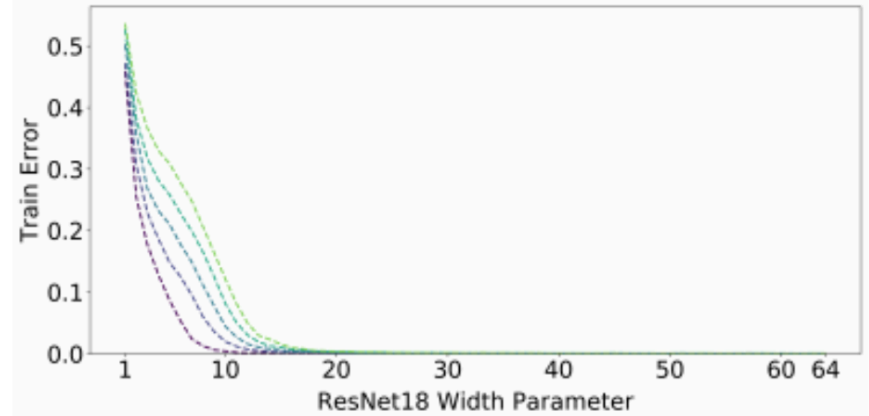
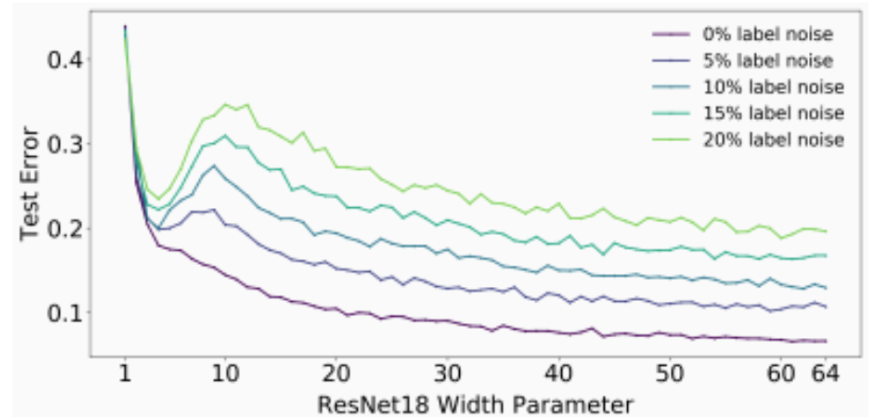
- There are cases where the model gets bigger, yet the (test!) loss goes down, sometimes even lower than in the classical “under-parameterized” regime.
- Complexity: number of parameters.

Double descent

Widespread phenomenon, across architectures (Nakkiran et al. '19):



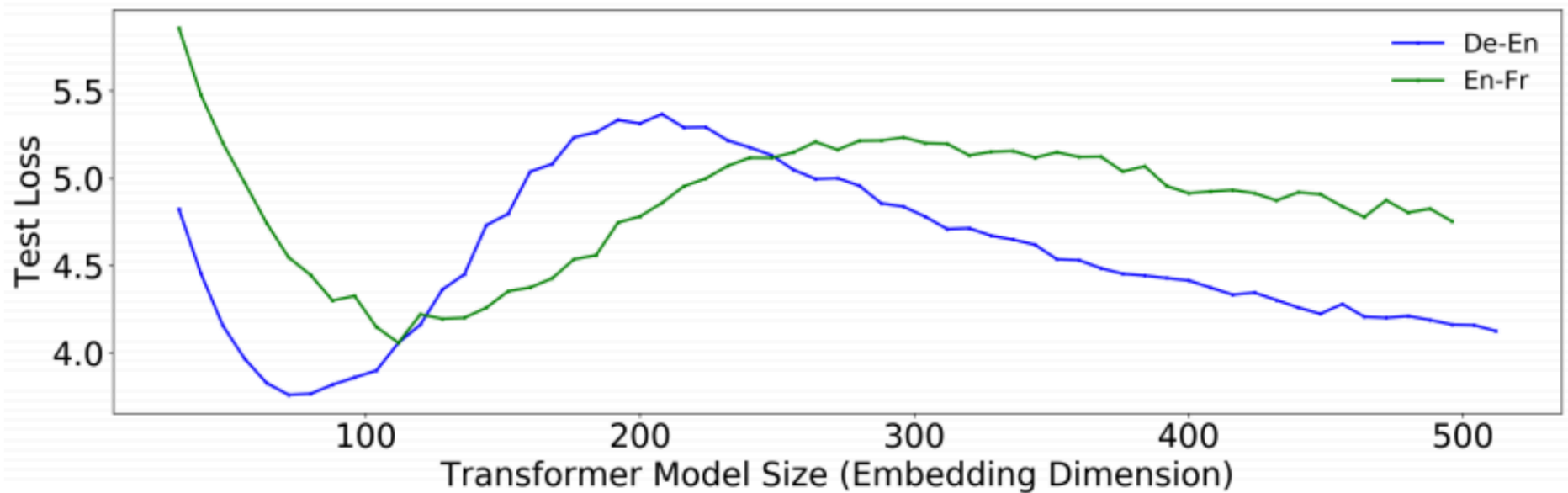
(a) **CIFAR-100.** There is a peak in test error even with no label noise.



(b) **CIFAR-10.** There is a “plateau” in test error around the interpolation point with no label noise, which develops into a peak for added label noise.

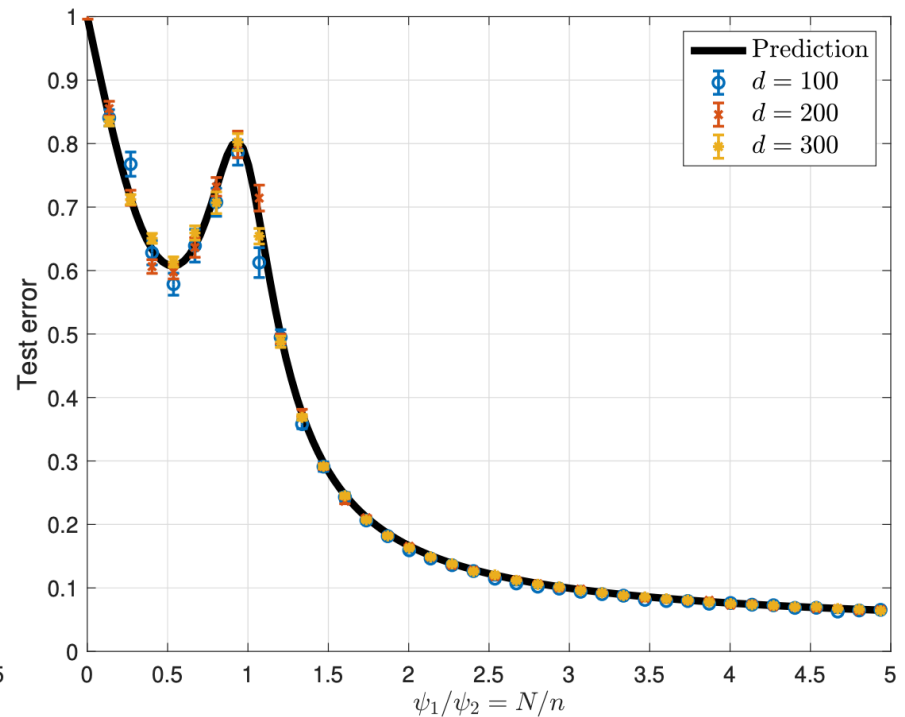
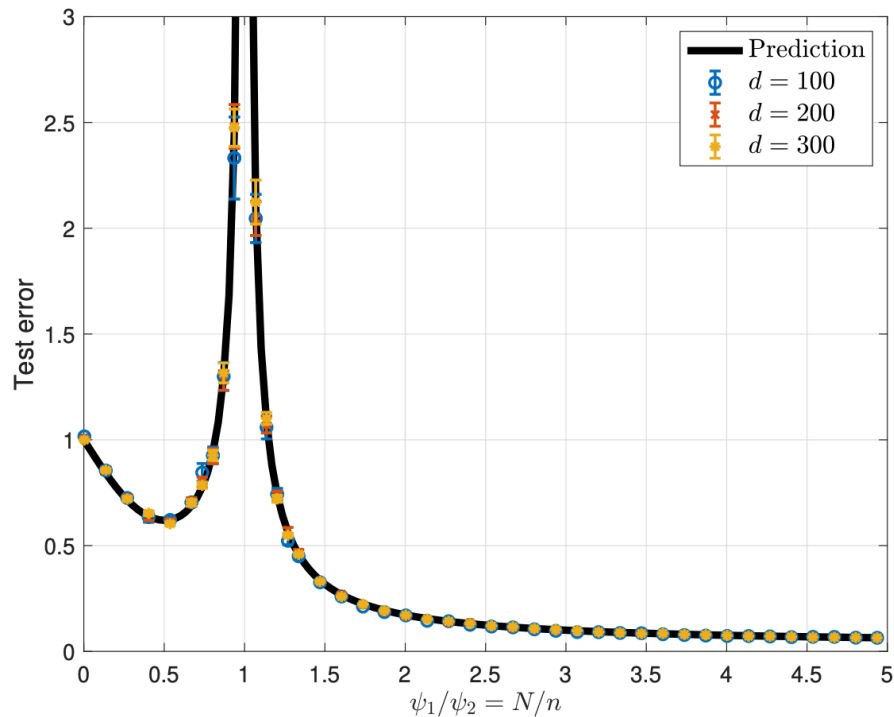
Double descent

Widespread phenomenon, across architectures (Nakkiran et al. '19):



Double descent

Widespread phenomenon, also in kernels (can be formally proved in some concrete settings [Mei and Montanari '20]), random forests, etc.



Double descent

Also in other quantities such as train time, dataset, etc (Nakkiran et al. '19):

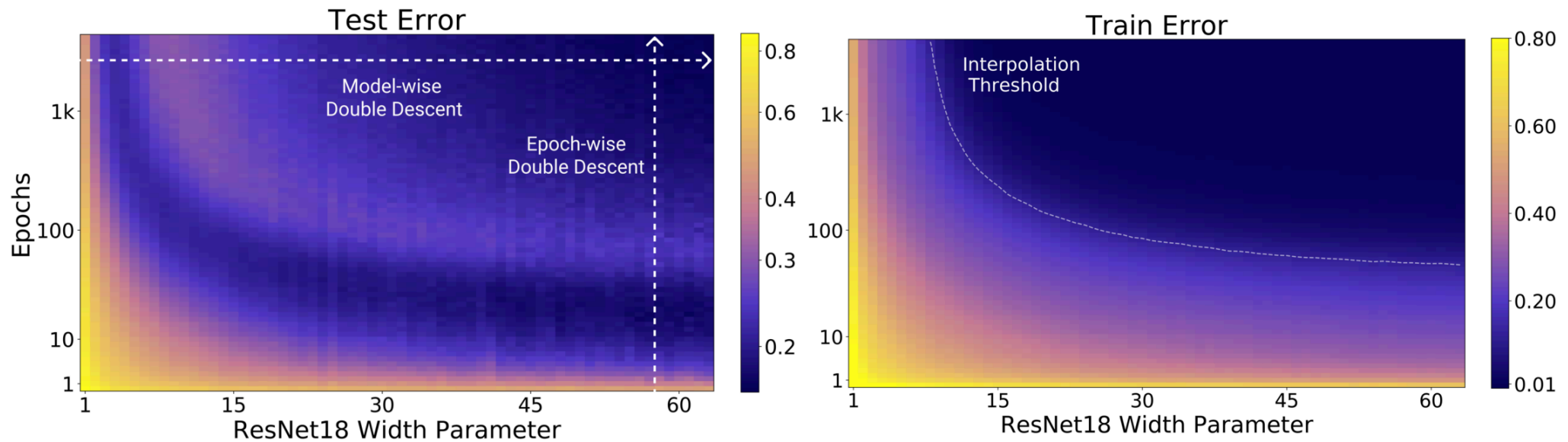
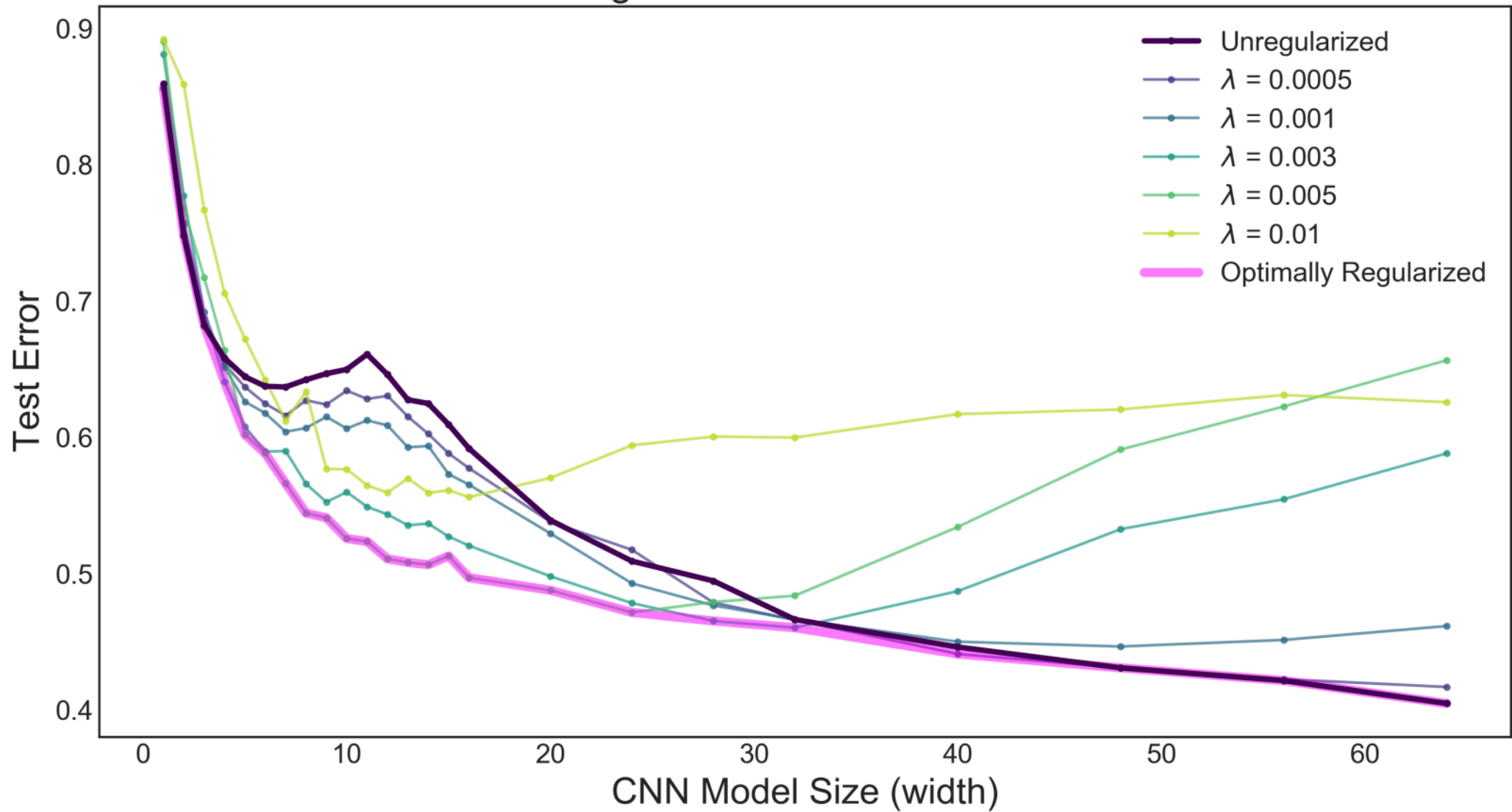


Figure 2: **Left:** Test error as a function of model size and train epochs. The horizontal line corresponds to model-wise double descent—varying model size while training for as long as possible. The vertical line corresponds to epoch-wise double descent, with test error undergoing double-descent as train time increases. **Right** Train error of the corresponding models. All models are Resnet18s trained on CIFAR-10 with 15% label noise, data-augmentation, and Adam for up to 4K epochs.

Double descent

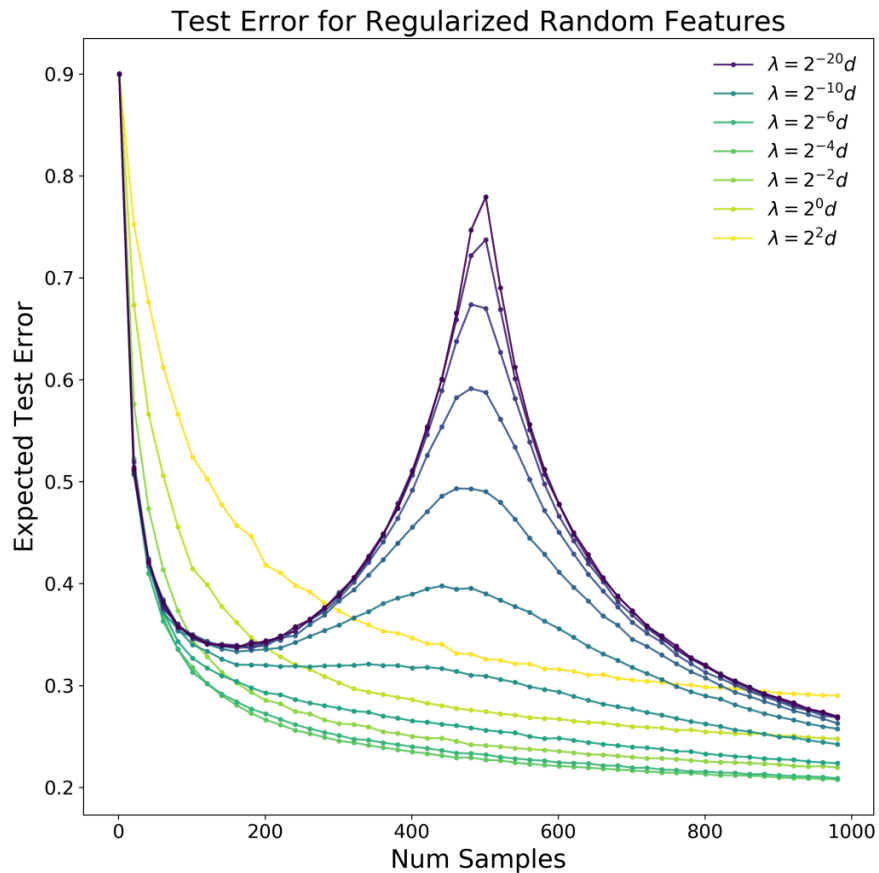
Optimal regularization can mitigate double descent [Nakkiran et al. '21]:

Effect of Regularization: CNNs on CIFAR-100

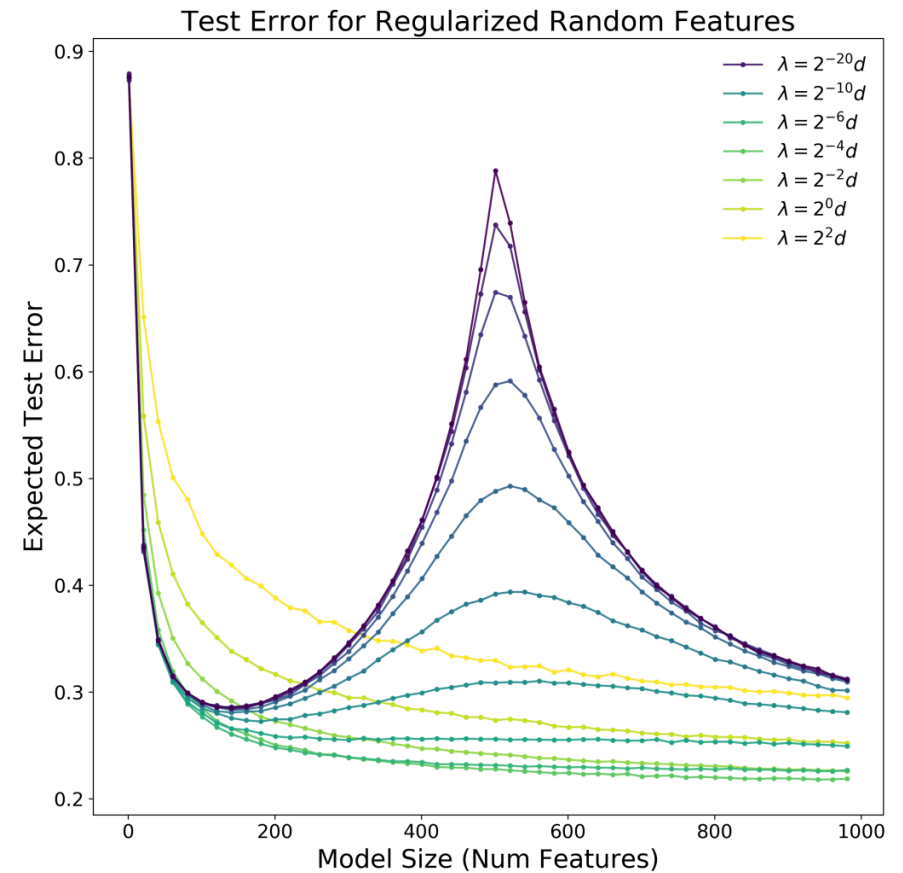


Double descent

Optimal regularization can mitigate double descent [Nakkiran et al. '21]:



a) Test Classification Error vs. Number of Training Samples.



(b) Test Classification Error vs. Model Size (Number of Random Features).

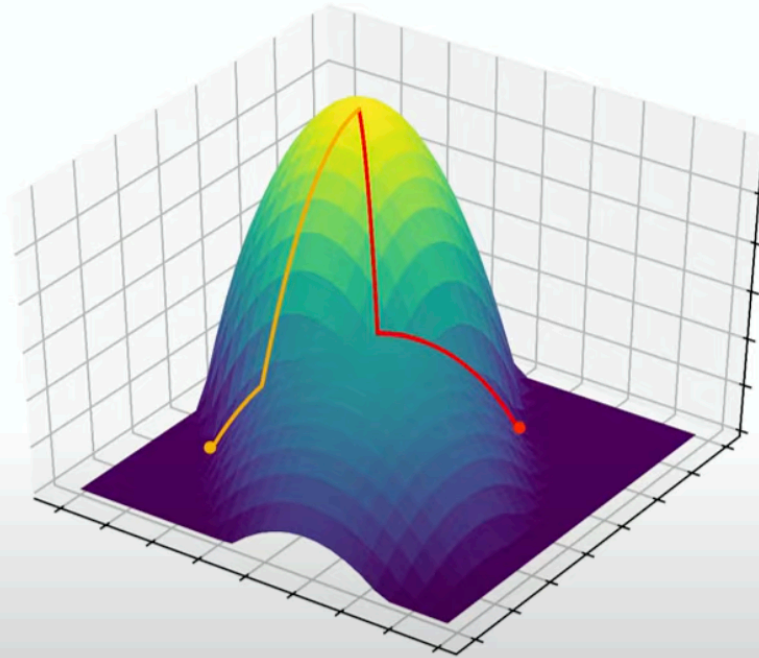
Implicit Regularization

Different optimization algorithm

→ Different bias in optimum reached

→ Different Inductive bias

→ Different generalization properties



Implicit Bias

Margin:

- Linear predictors:
 - Gradient descent, mirror descent, natural gradient descent, steepest descent, etc maximize margins with respect to different norms.
- Non-linear:
 - Gradient descent maximizes margin for homogeneous neural networks.
 - Low-rank matrix sensing: gradient descent finds a low-rank solution.