# Generalization Theory for Deep Learning

# Rademacher Complexity

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

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

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

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 $\ell$, we have

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

$$\underbrace{\qquad\qquad\qquad}_{\text{training error}}$$

$$\hat{R}_n \leq O(VC(\mathcal{F}))$$

and

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

# Norm-based Rademacher complexity bound

$\sigma: \text{ReLU}, \rho = 1$

**Theorem:** If the activation function is $\sigma$ is $\rho$-Lipschitz. Let

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

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

$\sigma(\S_y) \quad \to c_t \|_1 \text{ of } f$

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

$X \in \mathbb{R}^d, \quad W_1 \in \mathbb{R}^{m \times d}, \quad W_2, \dots, W_H \in \mathbb{R}^{m \times m}, \quad W_{H+1} \in \mathbb{R}^{m}$

$\|W_1^T\|_{1,\infty} = \max_{j=1,\dots,m} \|W_1^T(:,j)\|_1$

$= \max_{j=1,\dots,m} \sum_{j=1}^{d} |W_1^T(j,i)| \qquad W_1^T \, d \,\Box^m$

$\|W_H^T\|_{1,\infty} = \max_{j=1,\dots,m} \|W_H^T(:,i)\|_1 \qquad X^T \, n\,\Box^d$

$\|X^T\|_{2,\infty} = \max_{j=1,\dots,d} \|X^T(:,j)\|_2$

$= \max_{j=1,\dots,d} \|(X_1(j), X_2(j), \dots, X_n(j))\|_2$

# Massart Lemma

**Lemma:** Let $V$ be a set of vectors $\subset \mathbb{R}^d$, we have
$$R_V := \mathbb{E}_{\epsilon \sim \mathsf{Unif}\{1,-1\}^d} \max_{v \in V} \langle \epsilon, v \rangle \leq \sup_{v \in V} \|v\|_2 \sqrt{2 \ln d}.$$

- $\epsilon \sim \mathsf{Unif} \{1,-1\}^d$

- Fix $u \in V$

  Define $X_{u,i} = \epsilon_i \cdot u_i$

  $X_u = \sum_i \epsilon_i \cdot u_i$

  $(X_{u,i})^2 = u_i^2$

$\Rightarrow X_{u,i}$ is $u_i^2$-sub Gaussian

$\Rightarrow \|X_u\|_2$ is $\|u\|_2^2$-sub Gaussian

$\Rightarrow \forall u, \quad X_u$ is $\sup_{u \in V} \|u\|_2^2$ sub Gaussian

$\quad \mathbb{E}_\epsilon \max_{u \in V} \langle \epsilon, u \rangle \leq \sup_{u \in V} \|u\|_2 \cdot \sqrt{2 \ln d}$

<span style="color:red">Sub Gaussian expectation inequality</span>

# Properties of Rademacher Complexity

Let $V$ be a set of vectors $\subset \mathbb{R}^d$:

1. $R_{\text{Conv}(V)} = R_V$,

2. $R_V = R_{-V}$.

3. Let $V_1, \ldots, V_m$ be $m$ set of vectors such that for anyway $\epsilon \in \{-1, +1\}^d$, $\sup_{v \in V_i} \langle v, \epsilon \rangle \geq 0$ (e.g., $0 \in V_i$), then we have

$$R_{\cup_{i=1}^m V_i} \leq \sum_{i=1}^m R_{V_i}.$$

$u \in u$

$\Sigma \lambda i = 1$

$u \in \text{Conv}(V)$ iff $u = \sum_{i=1}^{u} \lambda_i u_i$ / $\lambda_i \geq 0$

# Proof of $(1,\infty)$ norm-based bound

Layer by layer induction

Let $\mathcal{F}_h = \{ x \longmapsto \sigma(W_h \sigma(W_{h-1} \cdots \sigma(W_1 x) \cdots) \in \mathbb{R}^{d_h}$

$$\|W_{h'}^T\|_{1,\infty} \leq B, \quad h' = 1, \ldots, h \}$$

$$\mathcal{F}_0 = \{ x \longmapsto x \}$$

Induction: $\mathrm{Rad}(\mathcal{F}_h|_S) \leq \|X^T\|_{2,\infty} (2\rho B)^h \sqrt{2 \ln d}$

take $h = l+1 \implies$ theorem

(1) $h = 0$ ✓

$\mathcal{F}_0 = \{ x \longmapsto (X(1), \ldots, X(d) : d \text{ functions} \}$

Massart lemma

$$\mathrm{Rad}(\mathcal{F}_0|S) \leq \left( \max_{j \in \{1, \ldots, d\}} \|(X_1(j), \ldots, X_n(j))\|_2 \right) \cdot \sqrt{2 \ln d}$$

$$= \|X^T\|_{2,\infty} \cdot \sqrt{2 \ln d} = \|X^T\|_{2,\infty} \cdot (2\rho B)^0 \sqrt{2 \ln d}$$

# Proof of $(1, \infty)$ norm-based bound

(2) Induction step, assume hypothesis is true till layer $h$

$$\mathcal{F}_{h+1}|_S = \{ x \mapsto \sigma(W_{h+1} \, g(x)), \, g \in \mathcal{F}_h, \, \|W_{h+1}\|_{1,\infty} \leq B \}$$

$$x \mapsto \sigma(W_{h+1} \, g(x)) \in \mathbb{R}^m$$

$$[\sigma(W_{h+1} \, g(x))](i),$$
$$= \sigma(W_{h+1}(i, :) \, g(x))$$

$$j = 1, \dots, m \quad \sum_i |V(i)_j|$$
$$\Longleftrightarrow = 1$$

$$V = \frac{W_{h+1}(i, :)}{\|W_{h+1}(i, :)\|_1}$$

$$\Longleftrightarrow \quad x \mapsto \sigma\left( \|W_{h+1}(i, :)\|_1 \cdot V^T g(x) \right), \quad \|V\|_1 = 1$$

$$\underbrace{\phantom{xxxxx}}_{\text{magnitude}} \quad \underbrace{\phantom{xxxxx}}_{\text{direction}}$$

$$= \sigma\left( \|W_{h+1}(i, :)\|_1 \cdot \sum_{j=1}^{m} V(j) \cdot [g(x)](j) \right)$$

$$= \sigma\left( \|W_{h+1}(i, :)\|_1 \cdot \sum_{j=1}^{m} |V(j)| \cdot \text{sgn}(V_j) [g(x)](j) \right)$$

$$\Longrightarrow \mathcal{F}_{h+1}|_S = \{ x \mapsto \sigma(\|W_{h+1}(i, :)\|_1 \cdot g(x)), \, g(x) \in \text{conv}(\mathcal{F}_h \cup (-\mathcal{F}_h)) \}$$

# Proof of $(1, \infty)$ norm-based bound

$$\text{Rad}\left(\mathcal{F}_{h+1}|_S\right)$$

$$\leq \ell B \; \text{Rad}\left(\text{Conv}\left(\left(\mathcal{F}_h|_S\right) \cup \left(-\mathcal{F}_h|_S\right)\right)\right)$$

$$= \ell B \; \text{Rad}\left(\mathcal{F}_h|_S \cup \left(-\mathcal{F}_h|_S\right)\right) \qquad \text{(property 1)}$$

$$\leq \ell B \left(\text{Rad}\left(\mathcal{F}_h|_S\right)\right) + \ell B \cdot \text{Rad}\left(-\mathcal{F}_h|_S\right) \qquad \text{(property 3)}$$

$$= 2\ell B \; \text{Rad}\left(\mathcal{F}_h|_S\right) \qquad \text{(property 2)}$$

$$\leq (2\ell B)^{h+1} \; \|X^T\|_{2,\infty} \; \sqrt{2 \ln d}$$

$$\square$$

# 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.
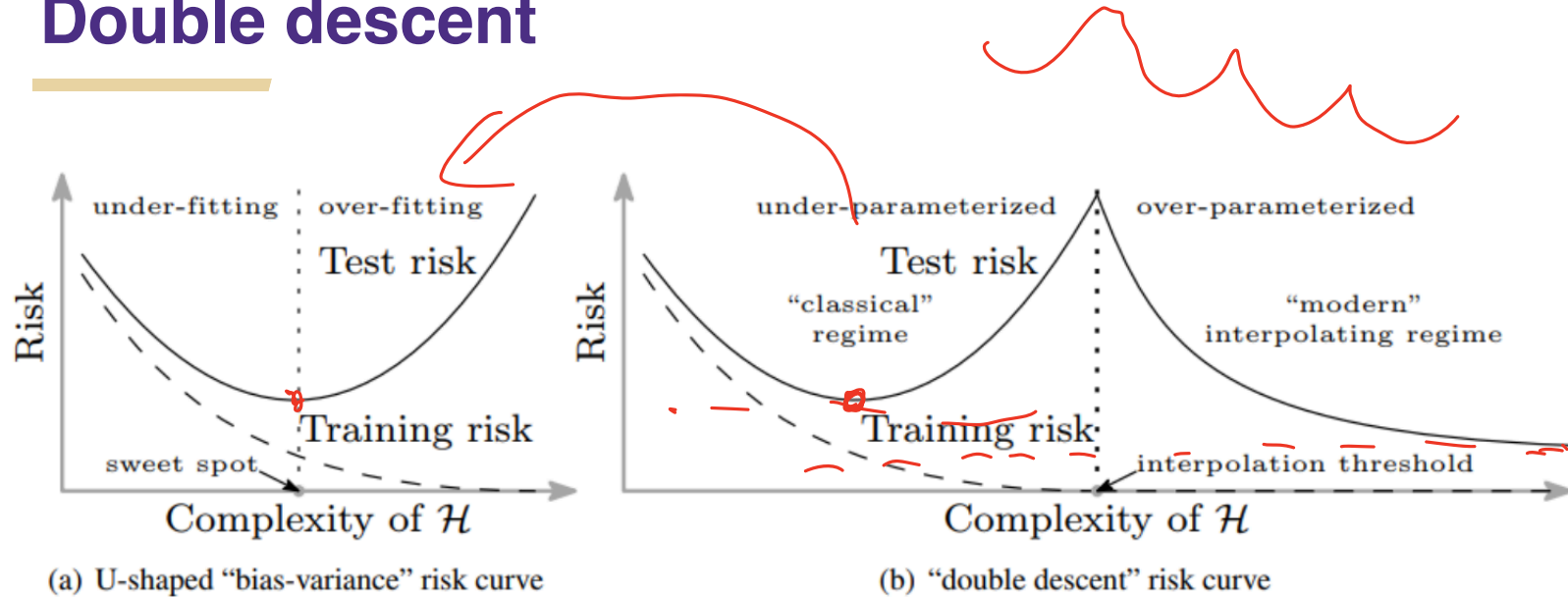


Image credits to Andrej Risteski

# 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.

  bias-variance decomposition

  K-nearest-neighbor

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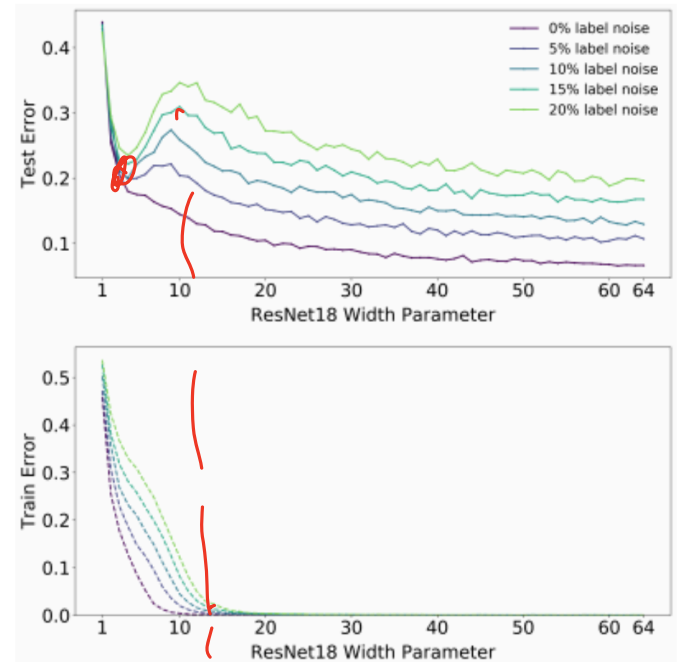  - Local Rademacher complexity: $1/n$.

# 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

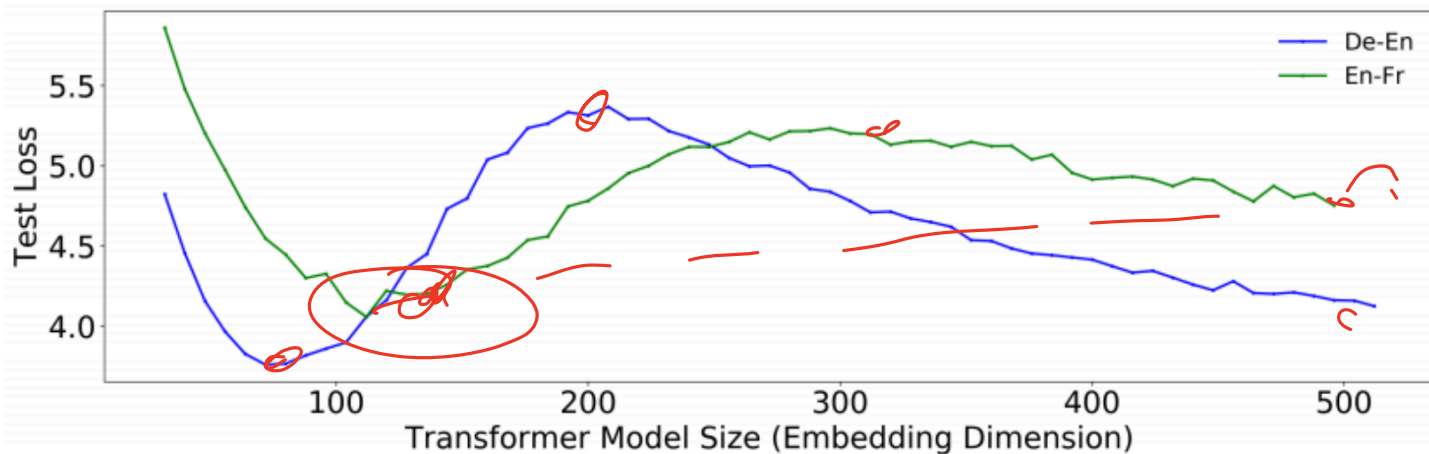Widespread phenomenon, across architectures (Nakkiran et al. '19):



(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.
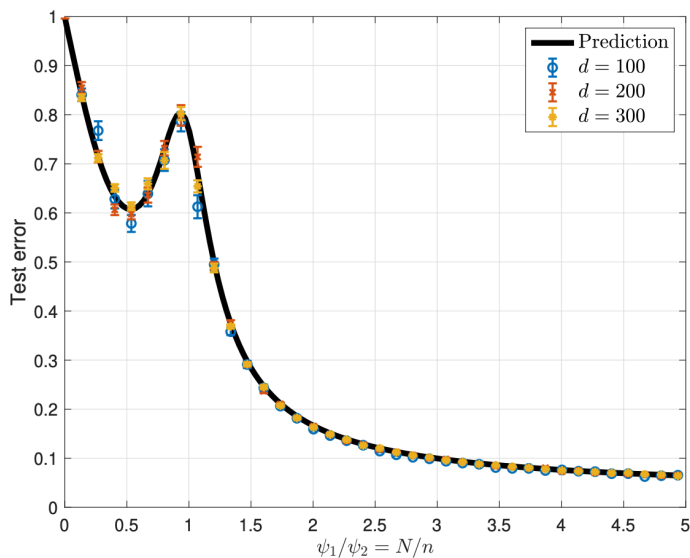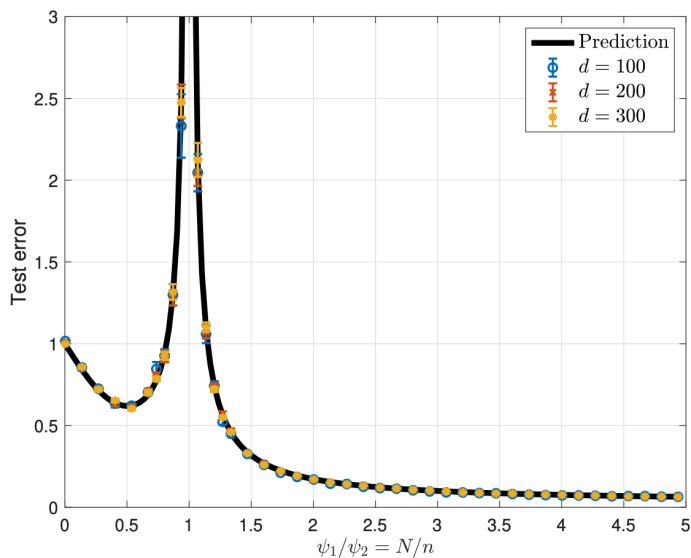
# Double descent

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

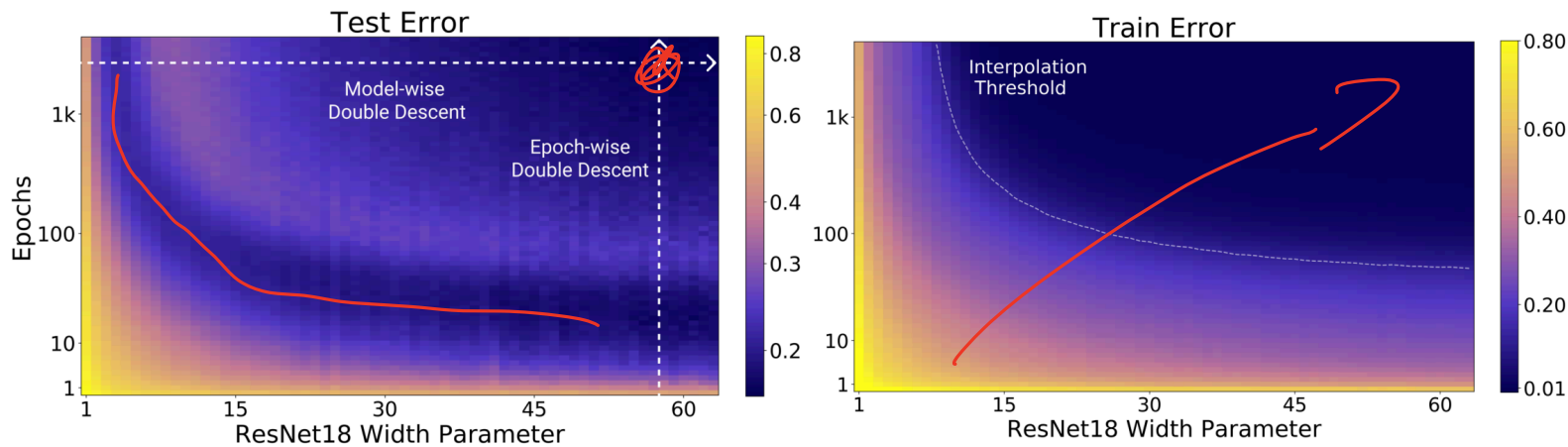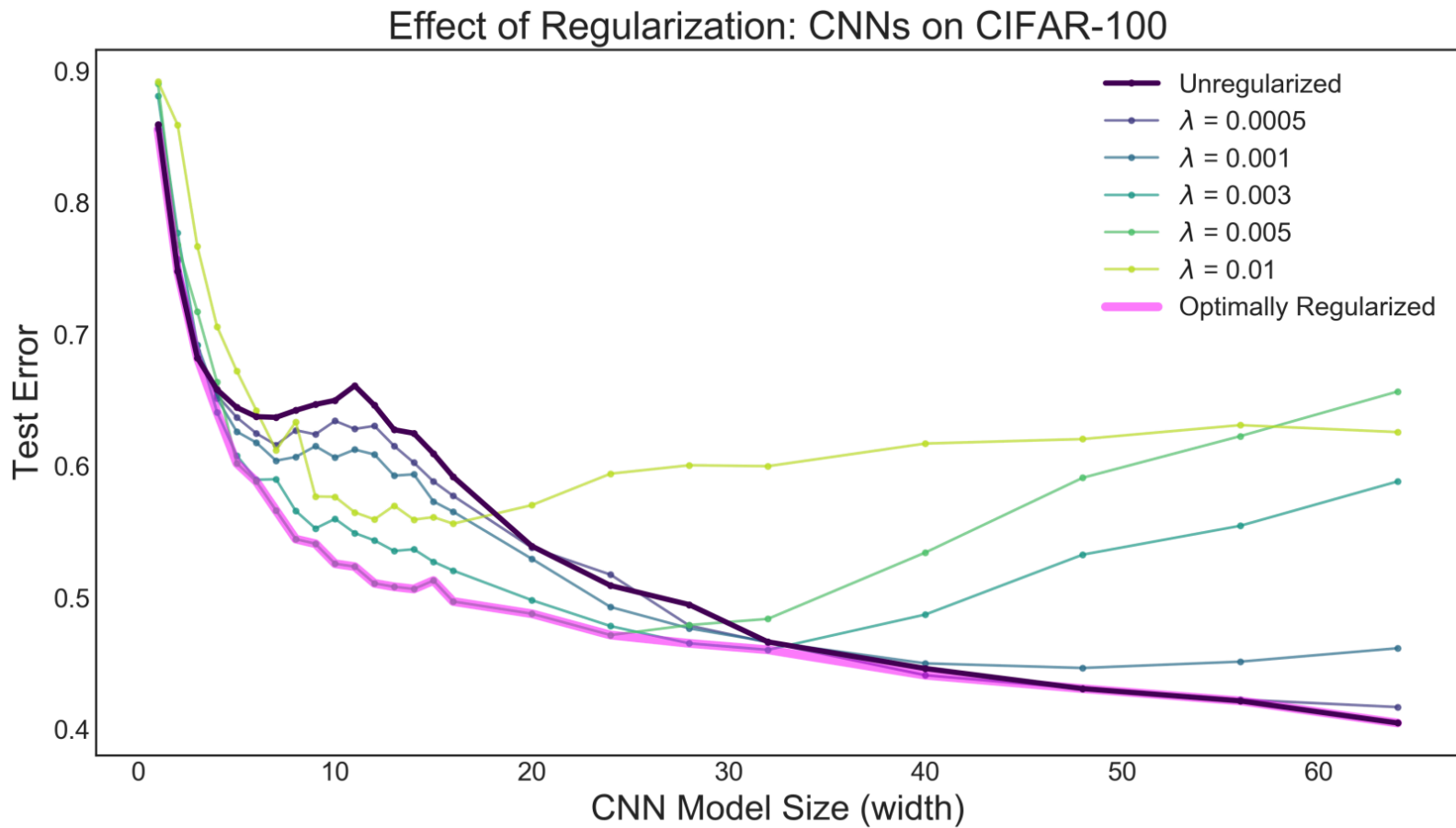Also in other quantities such as train time, dataset, etc (Nakkiran et al. '19):



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.
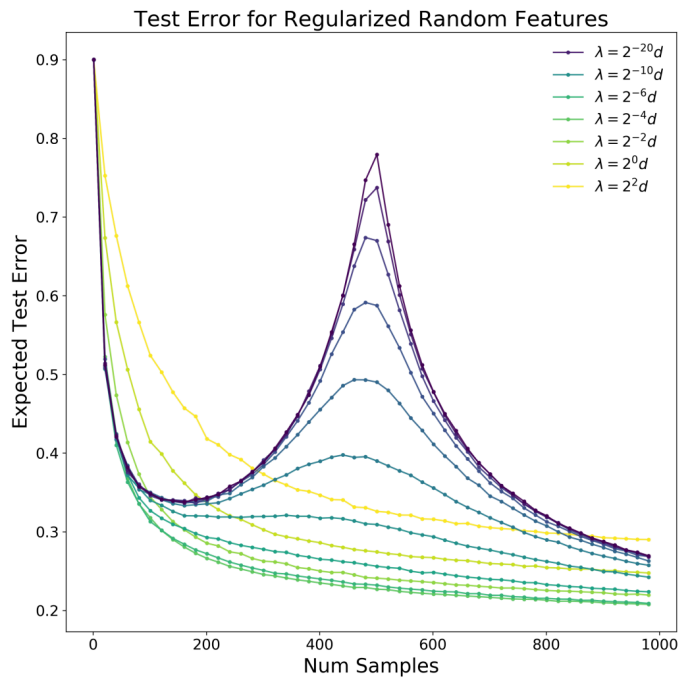
# Double descent

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



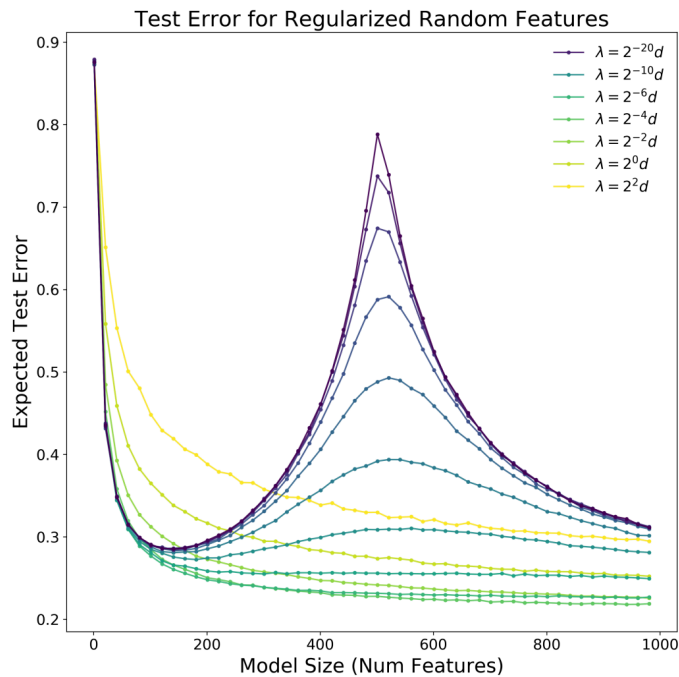Effect of Regularization: CNNs on CIFAR-100

# 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).
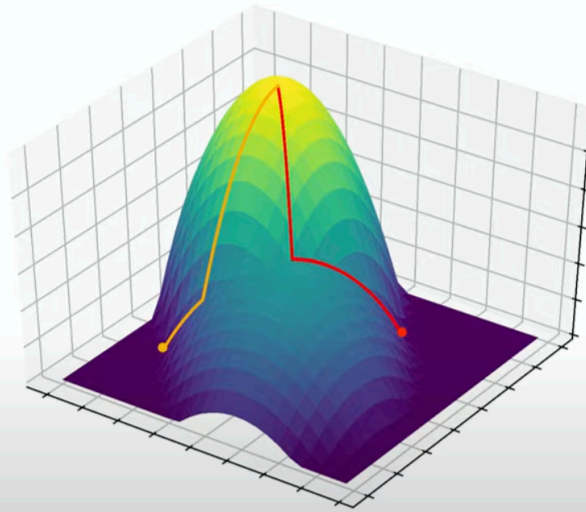
# Implicit Regularization

Different optimization algorithm
→ Different bias in optimum reached
→ Different Inductive bias
→ Different generalization properties



# params >> n
>/ global min
initialization, hyper-parameter

# Implicit Bias

**Margin:**

$linear$

$$\bar{u} = \arg\max_{\|w\|_2 \le 1}$$

$$\begin{aligned} &\min_i \\ &i = 1, \cdots y \end{aligned} \qquad y_i \langle w, x_i \rangle$$

$(1 - f(\cdot)) - homogeneous$

$$m_i(w) = y_i \cdot f(w, x_i)$$

$$\bar{u} = \frac{\min_i m_i}{\|w\|^{H+1}}$$

- 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.

$$M = UU^T \qquad \begin{aligned} &M: \text{ low-rank} \\ &U: \text{ full-rank} \end{aligned}$$

# 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 $\mathscr{F}$ such that $y = f(x)$ for some $f \in \mathscr{F}$:
  - no kernel is sample-efficient;
  - Exists a neural network that is sample-efficient.