



# Fast rates for noisy interpolation require rethinking the effects of inductive bias

June 9<sup>nd</sup> 2022, 1W-MINDS seminar

Fanny Yang, joint work with K. Donhauser, G. Wang, S. Stojanovic, N. Ruggeri

V Statistical Machine Learning group, CS department, ETH Zurich





#### Classical wisdom: Avoid fitting noise



Traditionally: want to avoid fitting noise perfectly for better (optimal) generalization.

# Double descent on neural networks

Classification using neural networks and first-order methods on CIFAR-10 with 15% label noise



After interpolation threshold, we have a second "descent" (double descent)

# Harmless interpolation on neural networks

Classification using neural networks and first-order methods on CIFAR-10 with 15% label noise



For large models, interpolation is not worse than regularization (harmless interpolation)

# Are these observations unique to neural networks?

Analogous surprising observations can be made for **linear models!** 

Here we run gradient descent on  $||y - Xw||_2^2$  with  $w_0 = 0$  for  $y = Xw^* + \xi$  with X,  $\xi$  standard Gaussians



# Interpolating models on two ends of analyzability

Neural network interpolators

- are found by using 1st order methods to minimize **non-convex losses**
- feature learning with overparameterization
  ~ e.g. width of hidden layers

Linear interpolators

• are interpolators that minimize



convex optimization problems

fixed *d* features with overparameterization
 ~ in terms of how much larger *d* >> *n*

complexity to analyze model

# High-level talk outline

- Observations for interpolating models
  - 。 "Second" descent
  - Harmless interpolation
- Explanation for the phenomena for linear interpolators
  - Previous intuition: Bias-variance trade-off by varying degree overparameterization
  - New complementary intuition: Bias-variance trade-off by varying inductive bias
  - Tight bounds show how moderate inductive bias can yield fast error rates

#### Previous work for min- $\ell_2$ -norm interpolators

Interpolators  $\hat{w} = \operatorname{argmin}_{w} ||w||_{2}$  s.t. y = Xw vs. Regularized estimator:  $||y - Xw||_{2}^{2} + \lambda ||w||_{2}^{2}$ Linear model  $y_{i} = \langle w^{*}, x_{i} \rangle + \xi_{i}$  with i.i.d.  $x_{i} \sim N(0, I)$ , some  $\xi_{i} \sim N(0, \sigma^{2})$ 



Previous bounds\* explain harmless interpolation & second descent: As  $\frac{d}{n}$  increases, variance decreases!

#### Previous work for min- $\ell_2$ -norm interpolators

Interpolators  $\hat{w} = \operatorname{argmin}_{w} ||w||_{2}$  s.t. y = Xw vs. Regularized estimator:  $||y - Xw||_{2}^{2} + \lambda ||w||_{2}^{2}$ Linear model  $y_{i} = \langle w^{*}, x_{i} \rangle + \xi_{i}$  with i.i.d.  $x_{i} \sim N(0, I)$ , some  $\xi_{i} \sim N(0, \sigma^{2})$ 



Are we happy? No, as opposed to NN, error overall is actually high for large  $\frac{d}{n}$  as the bias increases!

# What's missing? Structure...

Problem: Estimator has no "clue" where to search (all directions are equally valid)

• Line of work shows that  $\min \ell_2$ -interpolator can generalize well\*

but only for very specialized covariance  $\Sigma$  in practice  $\Sigma$  is fixed!

Question: What kind of interpolators can learn  $w^*$  well for large  $d \gg n$ ?

Classical intuition for  $d \gg n$ : good estimation only possible if

- we assume **simple structure of** *w*<sup>\*</sup> (such as sparsity) and
- the estimator has a strong matching inductive bias encouraging structural simplicity

Benefits of strong inductive bias (recap)

**Example for structural simplicity:** sparsity  $||w^*||_0 = s \ll d$ 

Estimators with weak (no) inductive bias: encouraging small  $||w||_2$  norm

**Matching strong inductive bias** : small  $||w||_0 / ||w||_1$  norm encouraging sparsity structure



# Old: Bias-variance trade-off via model complexity



But interpolators cannot attenuate noise-fitting by choosing an optimal  $\lambda$ !

# Our work: Bias-variance trade-off via inductive bias



Interpolators cannot attenuate by increasing *p* (decreasing structural bias)!

# Setting for presentation of our results (simplified)

- Function space: linear models  $f(x) = \langle w, x \rangle$  with  $x, w \in \mathbb{R}^d$
- Data model for *n* samples:  $y_i = \langle w^*, x_i \rangle + \xi_i$  with  $x_i \sim N(0, I)$  and noise  $\xi_i \sim N(0, \sigma^2)$

with sparse  $w^* = (1, 0, ..., 0)$  with unknown location (for simplicity of presentation)

- Degree of overparameterization:  $d \approx n^{\beta}$ ,  $\beta > 1$
- Minimum- $\ell_p$ -norm interpolators for  $p \in [1, 2]$ :  $\widehat{w} = \operatorname{argmin}_w ||w||_p$  s.t. y = Xw
- **Performance measure**: prediction error  $\mathbb{E}_{x \sim N(0,I)} (\langle x, \hat{w} w^* \rangle)^2 = ||\hat{w} w^*||^2$

# Recap for p = 2



For isotropic Gaussians,  $||\widehat{w} - w^*||^2 > c > 0$  for any  $\beta > 1$  ( $d \approx n^{\beta}$ ) even as  $n \to \infty$ 

# (Slow) rate for p = 1

Previous work for the i.i.d. noise case:

Theorem [WDY' 21](simplified) – Tight bounds for min- $\ell_1$ -norm interpolators

There exists a universal constant c > 0, s.t. whenever  $d = n^{\beta}$  with  $\beta > 1$ ,  $n \ge c$  w.h.p.

$$\left|\widehat{w} - w^{\star}\right| \Big|^{2} = \frac{\sigma^{2}}{\log\left(d/n\right)} + O\left(\frac{\sigma^{2}}{\log^{3/2}\left(d/n\right)}\right)$$

# (Slow) rate for p = 1

Theorem [WDY' 21](simplified) – Tight bounds for min- $\ell_1$ -norm interpolators

There exists a universal constant c > 0, s.t. whenever  $d \approx n^{\beta}$  with  $\beta > 1$ ,  $n \ge c$  w.h.p.

$$\left|\left|\widehat{w} - w^{\star}\right|\right|^{2} = \frac{\sigma^{2}}{\log\left(d/n\right)} + O\left(\frac{\sigma^{2}}{\log^{3/2}\left(d/n\right)}\right)$$



- This is a lower & upper bound for Gaussian X
- Experimentally, the bound is also tight beyond

#### Gaussian X, but hard to show!

Note: The same bound holds for classification

# (Slow) rate for p = 1

Theorem [WDY' 21](simplified) – Tight bounds for min- $\ell_1$ -norm interpolators

There exists a universal constant c > 0, s.t. whenever  $d \approx n^{\beta}$  with  $\beta > 1$ ,  $n \ge c$  w.h.p.

$$\left|\left|\widehat{w} - w^{*}\right|\right|^{2} = \frac{\sigma^{2}}{(\beta-1)\log n} + O\left(\frac{\sigma^{2}}{((\beta-1)\log n)^{3/2}}\right) \quad (\text{plugging in } d, n \text{ relation})$$

Second Descent after interpolation



Yes! Variance decreases, similar intuition as for p = 2



Harmless interpolation for large  $\beta$ 

No! Interpolator 
$$\Omega\left(\frac{1}{\log n}\right)$$
 vs. regularized  $O\left(\frac{s\log n}{n}\right)$ 

# Our work: Bias-variance trade-off via inductive bias



So far the extremes of very strong (p = 1) and no (p = 2) inductive bias perform badly

# Fast rates with $p \in (1,2)$

Theorem [DRSY' 22] (informal) – Upper & lower bounds for min- $\ell_p$ -norm interpolators

For  $d = n^{\beta}$  with  $1 < \beta \le \frac{p/2}{p-1'}$  and min- $\ell_p$ -norm interpolators with 1 and <math>n large enough,

we obtain with high probability, error rates of order  $\tilde{O}(n^{-\alpha})$  with  $\alpha$  as in graph below



better

- order-matching upper & lower bound
- for fixed  $\beta$ , some p > 1 close to 1 gets best rate
- for  $\beta \approx 2$ , rates close to  $\tilde{O}\left(\frac{1}{n}\right)$

Note: technique applies to classification (see paper) and allows extension to  $\Sigma \neq I$  and s-sparse w<sup>\*</sup>

# Fast rates with $p \in (1,2)$

Theorem [DRSY' 22] (informal) – Upper & lower bounds for min- $\ell_p$ -norm interpolators

For  $d \neq n^{\beta}$  with  $1 < \beta \leq \frac{p/2}{n-1}$  and min- $\ell_p$ -norm interpolators with 1 and <math>n large enough,

we obtain with high probability, error rates of order  $\tilde{O}(n^{-\alpha})$  with  $\alpha$  as in graph below



# Fast rates with $p \in (1,2)$ - caveat...

Theorem [DRSY' 22] (informal) – Upper & lower bounds for min- $\ell_p$ -norm interpolators

For  $d = n^{\beta}$  with  $1 < \beta \le \frac{p/2}{p-1'}$  and min- $\ell_p$ -norm interpolators with 1 and <math>n large enough,

we obtain with high probability, error rates of order  $\tilde{O}(n^{-\alpha})$  with  $\alpha$  as in graph below



Caveat:

"Large enough" actually requires

$$\frac{1}{\log \log d} \lesssim p - 1 \rightarrow \text{very large } \mathbf{d}$$

 $\Rightarrow$  cannot obtain best p for given  $\beta$ 

#### Our work: Bias-variance trade-off via inductive bias



# Bias-variance tradeoff for $p \in (1,2)$

For p = 1, variance and "sensitivity to noise" larger than for p = 2

 $\rightarrow$  increasing *d* vs. *n* does not regularize enough even though it has relatively small bias.



Trade-off between bias and variance for interpolators via strength of inductive bias!

#### Experimental results for classification

Experimental results: hard- $\ell_p$ -margin SVM for  $\sigma$ : proportion of random label flips



Isotropic Gaussians with  $d \sim 5000, n \sim 100$ 

Leukemia dataset with  $d \sim 7000, n \sim 70$ 

#### The tale of two "new" bias-variance trade-offs

#### **Previous intuition for interpolators:**

$$\widehat{w} = \operatorname{argmin}_{w} ||w||_{p} s.t.y = Xw$$

Bias-variance trade-off via overparameterization

#### **Our new intuition for interpolators**

$$\widehat{w} = \operatorname{argmin}_{w} ||w||_{p} s.t.y = Xw$$

Bias-variance tradeoff via inductive bias



decreasing effect of noise via increasing  $d_{eff}/n$  c

decreasing "strength of inductive bias" via increasing p

# Papers discussed in the talk



SML group: sml.inf.ethz.ch

- Wang\*, Donhauser\*, Yang "Tight bounds for minimum I1norm interpolation of noisy data", AISTATS '22
- Donhauser, Ruggeri, Stojanovic, Yang "Fast rates for noisy interpolation require rethinking the effects of inductive bias", ICML '22

