

# Generalization theory of Deep Learning

Yiqiao Zhong, Department of Statistics, UW Madison  
CS 762, Oct 20, 2022

# Some notations

- Training data  $S = (\mathbf{x}_i, y_i)_{i=1}^n$ , input  $x_i$  is  $d$ -dim vector,  $y_i$  label (or real value)
- Neural network (NN)  $f_{\theta}: R^d \rightarrow R^K$  (or  $R$ ), e.g., feed-forward NN

$$f_{\theta}(\mathbf{x}) = \mathbf{W}_L \sigma(\cdots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) + \mathbf{b}_L)$$

where  $\sigma$  is activation (e.g., ReLU);  $\theta$  contains all trainable parameters

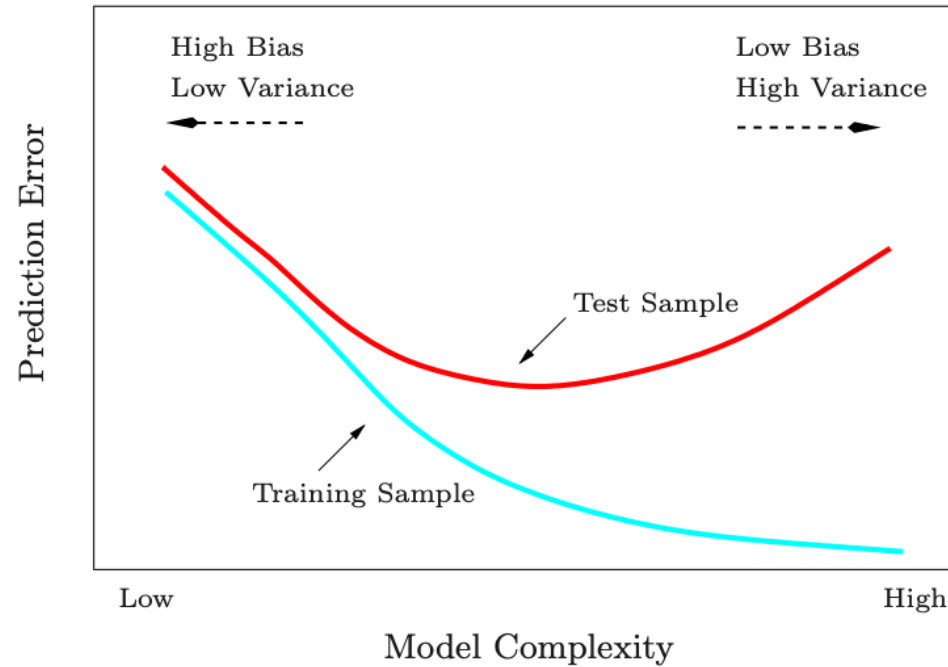
- Loss function  $L(y, f_{\theta}(\mathbf{x}))$ , Empirical Risk Minimization (ERM)

$$\hat{\theta} \in \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n L(y_i, f_{\theta}(\mathbf{x}_i)) =: R_n(\theta)$$

- Train loss  $R_n(\hat{\theta})$ , **train error**: ratio of misclassification on  $S$
- On test data  $T = (\mathbf{x}_i, y_i)_{i=1}^{n'}$ , evaluate  $\hat{\theta}$ , **test error** (sometimes generalization error): ratio of misclassification on  $T$ . Often  $n' = \infty$  in analysis
- Disclaimer: very incomplete references; check [Bartlett, Montanari, Rakhlin, Deep learning: a statistical view, 2021]

# The generalization puzzle

# Bias-variance tradeoff



**FIGURE 2.11.** Test and training error as a function of model complexity.

- Generally holds for many statistical models
- Classical solution to high-complexity models: regularize!

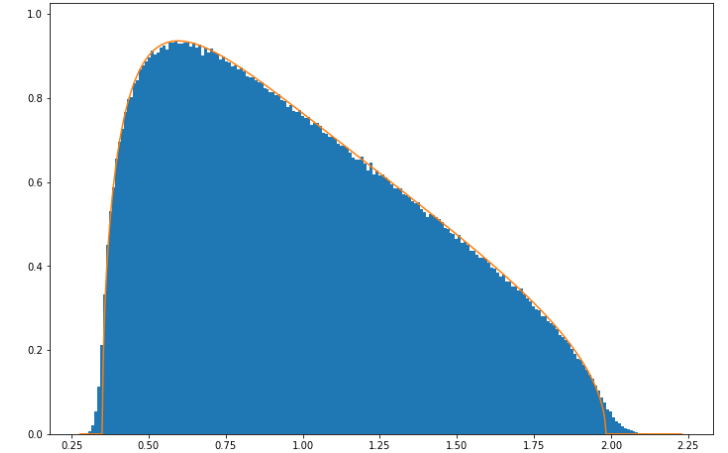
# Why and how regularizing high complexity model

- Consider linear ridge regression. Denote  $n \times d$  data matrix  $\mathbf{X}$ . Solve

$$\min_{\boldsymbol{\theta}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_2^2$$

- Using SVD :  $n^{-1/2}\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$

$$\hat{\boldsymbol{\theta}} = \left( \frac{1}{n} \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_d \right)^{-1} \frac{1}{n} \mathbf{X}^T \mathbf{y} = \frac{1}{\sqrt{n}} \sum_{j=1}^d \mathbf{u}_j \frac{\sigma_j}{\sigma_j^2 + \lambda} \mathbf{v}_j^T \mathbf{y}$$

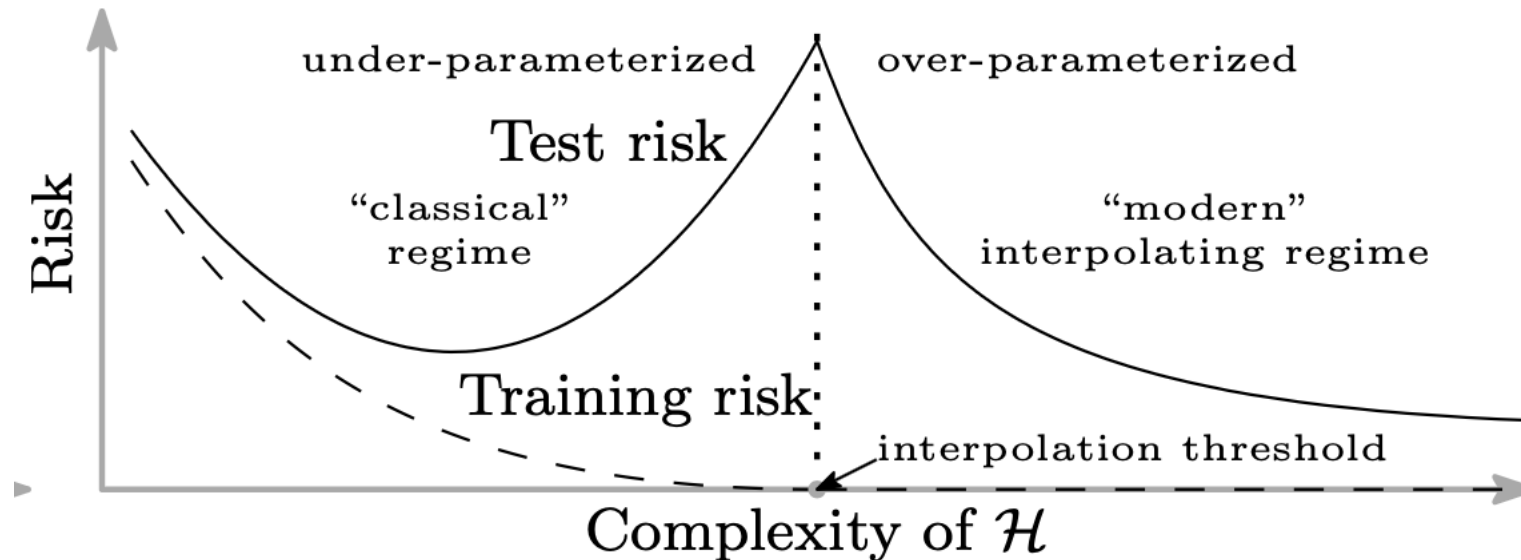


- Test error generally  $O(d/n)$  if  $d \ll n$
- Worse still, if  $d$  is close to  $n$ , **huge variance** in  $\hat{\boldsymbol{\theta}}$  without regularization. (MP law)
- Solution: need **large**  $\lambda$  if  $d$  is large.

Successful stories of regularization are everywhere:

- If signal is a sparse vector, use  $L_1$  regularization  $\|\boldsymbol{\theta}\|_1$ , called LASSO
- If signal is a low-rank matrix, use nuclear-norm regularization  $\|\boldsymbol{\theta}\|_*$

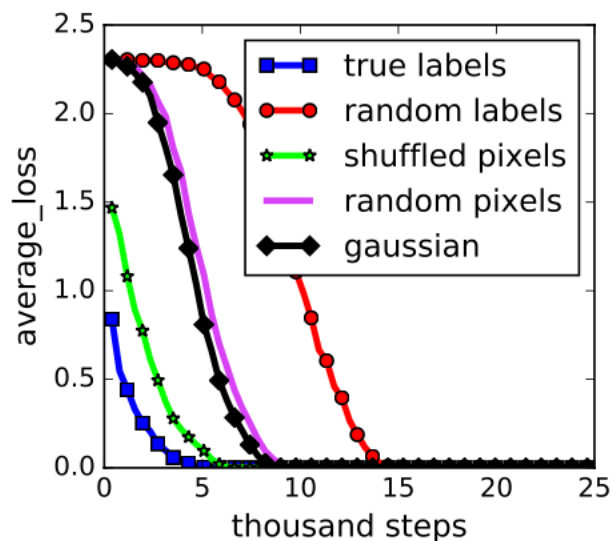
# But wait...double descent ?!



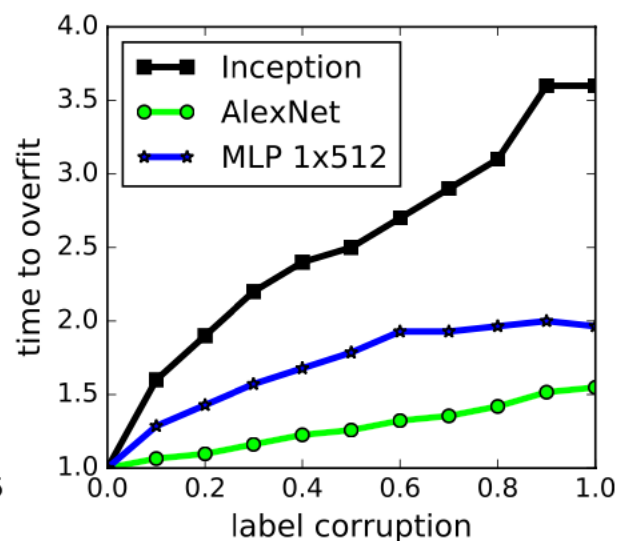
## New questions for ML/statistics:

1. When and why this happens?
2. When second descent better? Do we need regularization?
3. Lessons for architecture & algorithm design?

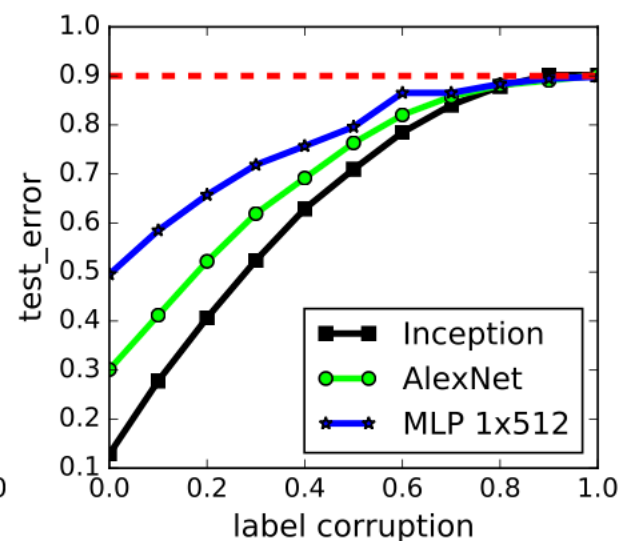
# Need understanding beyond interpolation



(a) learning curves



(b) convergence slowdown



(c) generalization error growth

Search for implicit bias



# Implicit bias

- Space of interpolating solutions (train error/loss is zero) may be large, but (stochastic) gradient descent (GD) converges to one with good generalization performance
- Proof-of-concept in overparametrized linear regression:

**Proposition** (HMRT, 2018). Initialize  $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ , and consider gradient descent on  $L_n(\boldsymbol{\beta}) = n^{-1}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ ,

$$\boldsymbol{\beta}^{(k)} = \boldsymbol{\beta}^{(k-1)} - \eta \nabla_{\boldsymbol{\beta}} L_n(\boldsymbol{\beta}), \quad k = 1, 2, 3, \dots,$$

where  $\eta > 0$  is sufficiently small. Then  $\lim_{k \rightarrow \infty} \boldsymbol{\beta}^{(k)} = \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}}$  is the minimum-norm interpolator:

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin} \left\{ \|\boldsymbol{\beta}\|_2 : \boldsymbol{\beta} \text{ minimizes } L_n(\boldsymbol{\beta}) \right\} = (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{y}.$$

Note that  $(\mathbf{X}^\top \mathbf{X})^+$  denotes the Moore-Penrose pseudoinverse. In particular, if  $\mathbf{X}\mathbf{X}^\top$  is invertible, then  $\hat{\boldsymbol{\beta}} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{y}$  is an interpolator.

# Implicit bias for classifying separable data

- Classification setting: for *linearly separable data* we can achieve zero train error using a linear classifier.

**Theorem** (Soudry, et al. 2018). Consider the logistic loss and any linearly separable data. From any initializer  $\beta^{(0)} \in \mathbb{R}^d$ , the gradient descent iterate  $\beta^{(k+1)} = \beta^{(k)} - \eta \nabla L_n(\beta^{(k)})$  satisfies

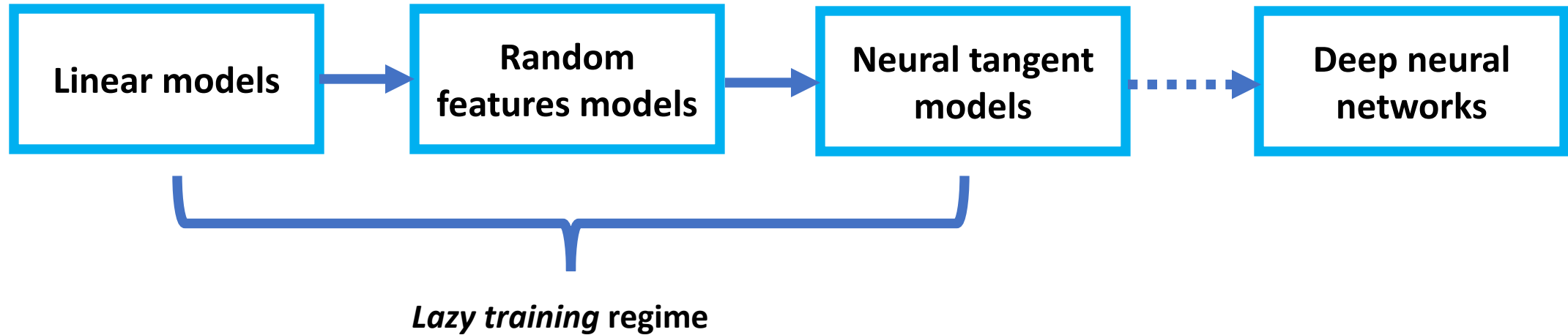
$$\beta^{(k)} = \hat{\beta} \log k + \Delta^{(k)}$$

where residual  $\|\Delta^{(k)}\|_2 = O(\log \log k)$  and  $\hat{\beta}$  is the max-margin solution

$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\beta\|_2^2 \quad \text{s.t. } \beta^\top \mathbf{x}_i \geq 1, \text{ for all } i = 1, 2, \dots, n.$$

- Gradient descent favors “small-norm” solution (at least in certain settings)
- Search for implicit bias: multiple linear deep network [Moroshko et al. 20], linear convolution network [GLSS18], one-hidden-layer ReLU network [NTS15], etc.
- Q: What is the generalization error of these solutions?

# The path to realism (or not?)



- In lazy training regime [OCB19], models are linear in parameters [HMRT18, MM19, MRSY19, MZ20]
- Test error can be calculated with idealized assumptions on data, rigorously justifying double descent

# Neural tangent model

- Key insight: when network width is infinite (or very large), the GD or SGD dynamics is given by (or approximated) by linearized local models---known as *neural tangent kernel* (NTK) models [JGH18, DZPS19, AZLS19, COB19]
- A simple example: one-hidden-layer NN:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{k=1}^N a_k \sigma(\mathbf{w}_k^\top \mathbf{x})$$

- Initialize from  $\boldsymbol{\theta}_0 = (\mathbf{a}_0, \mathbf{W}_0)$ , do **Taylor expansion**:

$$\begin{aligned} & \frac{1}{\varepsilon} f(\mathbf{x}; \mathbf{a}_0 + \varepsilon \mathbf{a}, \mathbf{W}_0 + \varepsilon \mathbf{W}) \\ & \approx \underbrace{\frac{1}{\varepsilon} f(\mathbf{x}; \mathbf{a}_0, \mathbf{W}_0)}_{\text{initialization}} + \underbrace{\sum_{k=1}^N a_k \sigma(\langle \mathbf{w}_{0,k}, \mathbf{x} \rangle)}_{\text{random features model}} + \underbrace{\sum_{k=1}^N a_{0,k} \langle \mathbf{w}_k, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_{0,k}, \mathbf{x} \rangle)}_{\text{neural tangent model}} \end{aligned}$$

# Why NTK makes sense?

Consider residuals at time  $t$ : let  $\mathbf{r}_t = \mathbf{y} - f(\mathbf{x}; \boldsymbol{\theta}_t) \in \mathbb{R}^n$ , then running gradient flow (GF) gives

$$\frac{d}{dt} \mathbf{r}_t = -\mathbf{K}_t \mathbf{r}_t, \quad \text{where } [\mathbf{K}_t]_{ij} = \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i; \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_j; \boldsymbol{\theta}_t) \rangle$$

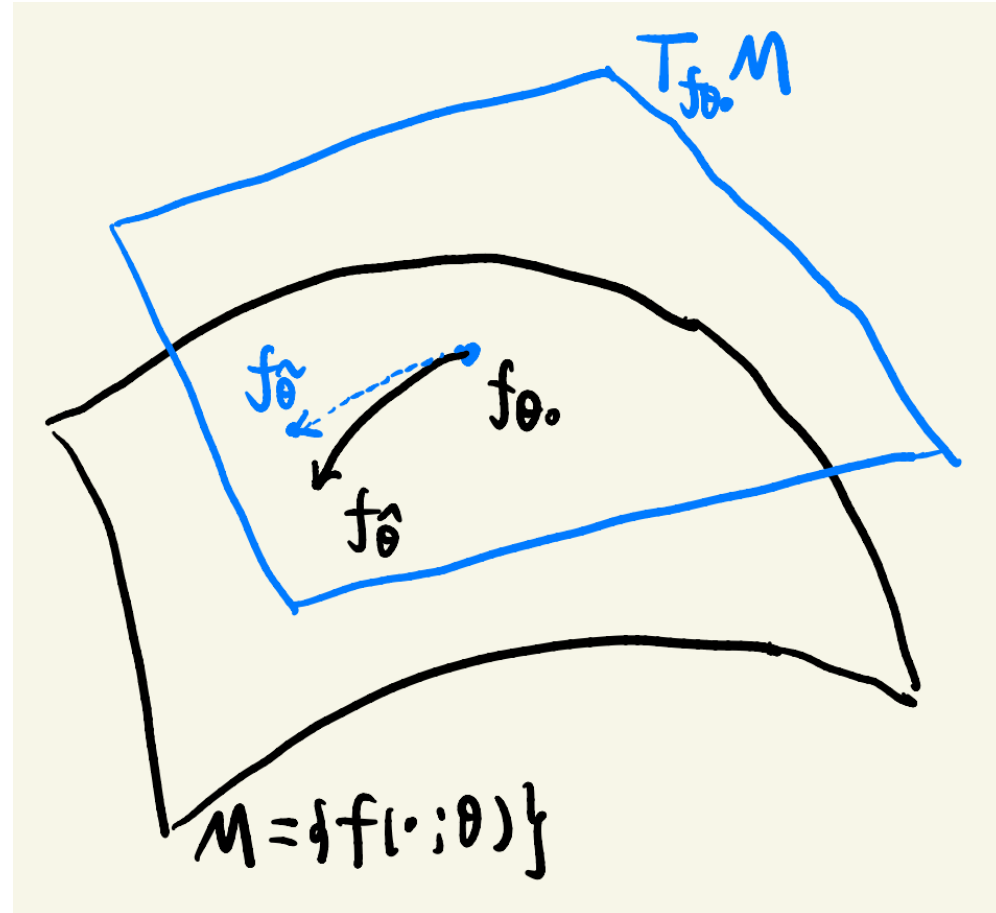
If  $\mathbf{K}_t \approx \mathbf{K}_0$  for all  $t$  (guaranteed when width is large), then

$$(i) \quad \frac{d}{dt} \|\mathbf{r}_t\|_2^2 = -2\mathbf{r}_t^\top \mathbf{K}_t \mathbf{r}_t \lesssim -2\lambda_{\min}(\mathbf{K}_0) \|\mathbf{r}_t\|_2^2$$

$\implies$  linear convergence rate

$$(ii) \quad \mathbf{r}_t \approx \tilde{\mathbf{r}}_t \quad \text{where } \tilde{\mathbf{r}}_t \text{ comes from GF}$$

$$\text{on the loss } \tilde{L}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\text{lin}}(\mathbf{x}_i; \boldsymbol{\theta}))^2$$



- Q: under NTK, what is the generalization error?

Insights from statistics

# Overparametrized linear models

- Consider the one-hidden-layer NTK model. We have  $NT$  features:

$$\nabla_{\mathbf{w}} f(\mathbf{x}; \boldsymbol{\theta}_0) = [a_1 \sigma'(\mathbf{x}^\top \mathbf{w}_{0,1}) \mathbf{x}^\top, \dots, a_N \sigma'(\mathbf{x}^\top \mathbf{w}_{0,N}) \mathbf{x}^\top] \in \mathbb{R}^{Nd}$$

- A useful simplification: NT features have complicated dependence, why not assume that we have  $\tilde{\mathbf{x}}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ , prediction function is  $\tilde{\mathbf{x}}^\top \hat{\boldsymbol{\theta}}$ . By abuse of notations, just write  $\mathbf{x}_i$ .

**Assumption.** Suppose  $\mathbf{z} = \boldsymbol{\Sigma}^{-1/2} \mathbf{x}$  is 1-subgaussian.

WLOG, assume  $\boldsymbol{\Sigma} = \text{diag}(\lambda_1, \dots, \lambda_d)$  where  $\lambda_1 \geq \dots \geq \lambda_d$ .

- Key insight:

$$\langle \mathbf{x}, \hat{\boldsymbol{\theta}} \rangle = \underbrace{\langle \mathbf{x}_{\leq k}, \hat{\boldsymbol{\theta}}_{\leq k} \rangle}_{\text{prediction part}} + \underbrace{\langle \mathbf{x}_{>k}, \hat{\boldsymbol{\theta}}_{>k} \rangle}_{\text{interpolation part}}$$

# Decomposing features

- Regression setting. Data  $\mathbf{X} = [\mathbf{X}_{\leq k}, \mathbf{X}_{> k}]$  of size  $n \times d$ , where  $d > n$

$$\mathbf{X}\mathbf{X}^\top = \mathbf{X}_{\leq k}\mathbf{X}_{\leq k}^\top + \mathbf{X}_{> k}\mathbf{X}_{> k}^\top$$

- A seemingly bold assumption:  $\mathbf{X}_{> k}\mathbf{X}_{> k}^\top \approx \gamma \mathbf{I}_n$
- Heuristic justification: features are divided into “important” ones ( $\leq k$ ) and “not important” ones ( $> k$ ); the latter is similar to pure noise

$$\hat{\boldsymbol{\theta}}_{\leq k} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^k} \|\mathbf{y} - \mathbf{X}_{\leq k}\boldsymbol{\theta}\|_2^2 + \gamma \|\boldsymbol{\theta}\|_2^2$$

- Equivalent to **ridge regression!**



# Implicit regularization

- Parameter  $\gamma$  controls the amount of regularization
- Turning heuristics into rigorous arguments. For general  $(\lambda_j)_{j \leq d}$ , define *effective rank*:

$$r_k(\boldsymbol{\Sigma}) = \frac{\sum_{i>k} \lambda_i}{\lambda_{k+1}}.$$

- Concentration results can show:

$$r_k(\boldsymbol{\Sigma}) \geq bn \implies \lambda_j(\mathbf{X}_{>k} \mathbf{X}_{>k}^\top) \in [\gamma/c, c\gamma] \text{ where } \gamma = \sum_{j>k} \lambda_j$$

- Find a sweet spot for  $k$  so that:  $\mathbf{x}_{\leq k}$  captures almost all information while  $\mathbf{x}_{>k}$  is similar to noise. Called *effective dimension*.

# A look at the theorems

**Theorem 4.4.** Fix  $\delta < 1/2$ . Under Assumption 4.2, suppose for some  $k$  the condition number of  $\mathbf{X}_{>k}\mathbf{X}_{>k}^\top$  is at most  $\kappa$  with probability at least  $1 - \delta$ . Then

$$\widehat{\text{VAR}} \lesssim \sigma_\xi^2 \kappa^2 \log\left(\frac{1}{\delta}\right) \left(\frac{k}{n} + \frac{n \sum_{i>k} \lambda_i^2}{(\sum_{i>k} \lambda_i)^2}\right) \quad (44)$$

with probability at least  $1 - 2\delta$ .

**Theorem 4.5.** Under the assumptions of Theorem 4.4, for  $n \gtrsim \log(1/\delta)$ , with probability at least  $1 - 2\delta$ ,

$$\widehat{\text{BIAS}}^2 \lesssim \kappa^4 \left[ \|\boldsymbol{\theta}_{\leq k}^*\|_{\boldsymbol{\Sigma}_{\leq k}^{-1}}^2 \left(\frac{\sum_{i>k} \lambda_i}{n}\right)^2 + \|\boldsymbol{\theta}_{>k}^*\|_{\boldsymbol{\Sigma}_{>k}}^2 \right]. \quad (48)$$

- Upper bounds tight up to constants [Tsigler, Bartlett, 2020]
- Bias and variance vanish under suitable decay of eigenvalues [TB20], empirically checked [WHS22]

# Is linear model naïve?

- Consider the NT features:

$$\underbrace{\sigma'(\langle \mathbf{x}, \mathbf{w}_k \rangle) \mathbf{x}}_{\text{NT feature}} = \underbrace{\sigma'_{\leq \ell}(\langle \mathbf{x}, \mathbf{w}_k \rangle) \mathbf{x}}_{\text{effective fitting}} + \underbrace{\sigma'_{> \ell}(\langle \mathbf{x}, \mathbf{w}_k \rangle) \mathbf{x}}_{\text{noise}}$$

- The spirit is the same. Stacking NT features into  $n \times (Nd)$  matrix  $\Phi$ . Assume isotropic data  $\mathbf{x}_i \sim N(\mathbf{0}, \mathbf{I}_d)$ .

$$\begin{aligned} \Phi \Phi^\top &\approx \Phi_{\leq k} \Phi_{\leq k}^\top + \Phi_{> k} \Phi_{> k}^\top && \text{(cross term negligible)} \\ &\approx \Phi_{\leq k} \Phi_{\leq k}^\top + \|\sigma'_{> k}\|_{L^2}^2 \cdot \mathbf{I}_n && \text{(high-degree term concentrates)} \end{aligned}$$

- *Self-induced regularization*: **nonlinearity** of activation helps!

# A general generalization result for 2-layer NTK

- [Montanari, Zhong, 2020] Suppose  $d^k \ll n \ll d^{k+1}$ , isotropic input data. general target function  $f_* \in L_2(S^{d-1})$ . As long as network width  $N$  satisfies  $Nd \gg n$  (overparametrization), then with high probability,

$$\begin{aligned} R_{\text{NT}}(f; \lambda) &= R_{\text{KRR}}(f; \lambda) + O\left(\sqrt{\frac{n}{Nd}}\right) \\ &= R_{\text{PRR}}(f; \lambda + \|\sigma'_{>k}\|_{L^2}^2) + O\left(\sqrt{\frac{n}{Nd}}\right) \end{aligned}$$

- Generalization via low-degree component, interpolation via high-degree component
- Regularization increased due to high-degree part of activation

# Beyond Lazy Training

# Limitation of Lazy training

- Success of deep learning depends on *learning data representation*. More complicated than random features models or variants.
- Want NNs to move moderately away from initialization.
- Random features models restricted, having trouble learning single neuron function [MBM17].
- Nevertheless, NTK may be advantageous for small-sample datasets [ADLS+19]

# Mean-field perspective

- Viewing parameters as a *probability distribution* [MMN18, CB18]

$$f(\mathbf{x}; \boldsymbol{\theta}) = \int a \sigma(\mathbf{w}^\top \mathbf{x}) \hat{\rho}_N(da, d\mathbf{w}), \quad \text{where } \hat{\rho}_N \text{ is empirical measure on } \mathbb{R}^{d+1}$$

- Under a nonstandard initialization scaling  $\text{var}(w_k) \sim O\left(\frac{1}{n^2}\right)$ , continuous-time SGD  $\approx$  Gradient flow on probability measure, which is determined by a PDE.
- Advantage: capable of learning more functions
- Disadvantages: weak theory. Requires very large width (likely exponential in  $d$ ), requires very large sample size (in general, exponential in  $d$ ); the latter can be improved to polynomial dependence by adding noise [WLLM20]

# Feature learning with GD

- Suppose the target function  $f_*(\mathbf{x}) = g(\mathbf{U}\mathbf{x})$  where  $\mathbf{U}$  is of size  $d \times r$  with  $d \gg r$ . Assume  $g$  is of polynomial of degree  $p$ .
- NTK cannot learn the unknown subspace  $\mathbf{U}$ , thus requiring a much larger sample size  $O(d^p)$
- Assuming non-degeneracy condition of Hessian of  $f_*$ , one-step GD on the squared loss using one-hidden-layer NN reveals information about  $\mathbf{U}$ , which results in improved sample complexity  $O(d^2)$ ; see [DLS22]



# Other approaches

- Classical tools in learning theory such as VC dimension insufficient because dimension is too large [BMM99]
- Bounding Rademacher using weight matrix norms [BFT17]
- Finding other good complexity measures by taking into account initialization [NLBLS18], algorithms, etc.

**Theorem 2.** For any  $h \geq 2$ ,  $\gamma > 0$ ,  $\delta \in (0, 1)$  and  $\mathbf{U}^0 \in \mathbb{R}^{h \times d}$ , with probability  $1 - \delta$  over the choice of the training set  $\mathcal{S} = \{\mathbf{x}_i\}_{i=1}^m \subset \mathbb{R}^d$ , for any function  $f(\mathbf{x}) = \mathbf{V}[\mathbf{U}\mathbf{x}]_+$  such that  $\mathbf{V} \in \mathbb{R}^{c \times h}$  and  $\mathbf{U} \in \mathbb{R}^{h \times d}$ , the generalization error is bounded as follows:

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

Emerging phenomena, and new hope?

# Self-supervised learning

- Representation using supervised learning  $f(\mathbf{x}; \boldsymbol{\theta})$ . Q: label intensive? How to transfer?
- With no (or very few) label information, NNs can learn good embedding, e.g., SimCLR [CKNH20]
- Clear cluster structure & meaningful learned features
  
- Self-supervised learning or unsupervised learning may be a bridge to understanding generalization



# Visualizing learned features

## Supervised Features

(of adversarial trained Wide-ResNet)



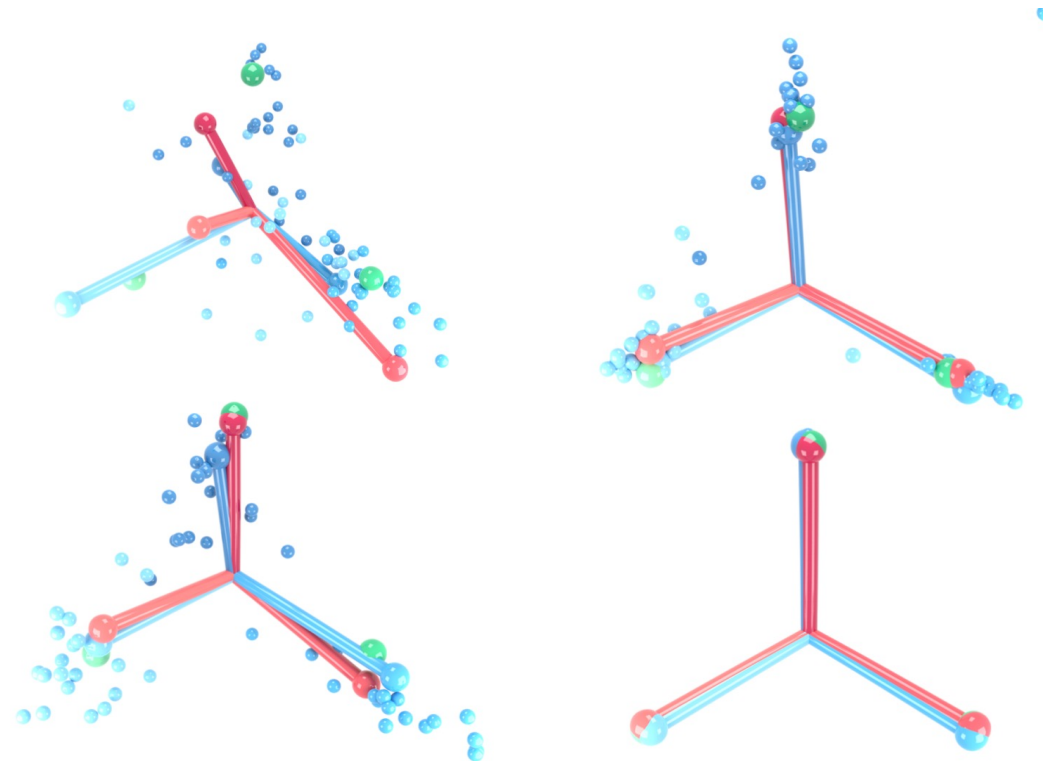
## Contrastive Features

(of adversarial-contrast trained Wide-ResNet)



# Neural collapse

- $h = h_\theta(x)$  is last-layer activations, where  $h_\theta: R^d \rightarrow R^p$ ,  $K$  classes
- Classifier:  $Wh_\theta(x) + b$
- decomposing covariance: between-class + within-class:
- $\Sigma_T = \Sigma_B + \Sigma_W$



# Clear phenomenon, clean math relations

(NC1) Variability collapse:  $\Sigma_W \rightarrow 0$

(NC2) Convergence to Simplex ETF:

$$\begin{aligned} & \left| \|\boldsymbol{\mu}_c - \boldsymbol{\mu}_G\|_2 - \|\boldsymbol{\mu}_{c'} - \boldsymbol{\mu}_G\|_2 \right| \rightarrow 0 \quad \forall c, c' \\ & \langle \tilde{\boldsymbol{\mu}}_c, \tilde{\boldsymbol{\mu}}_{c'} \rangle \rightarrow \frac{C}{C-1} \delta_{c,c'} - \frac{1}{C-1} \quad \forall c, c'. \end{aligned}$$

(NC3) Convergence to self-duality:

$$\left\| \frac{\mathbf{W}^\top}{\|\mathbf{W}\|_F} - \frac{\dot{\mathbf{M}}}{\|\dot{\mathbf{M}}\|_F} \right\|_F \rightarrow 0$$

(NC4): Simplification to NCC:

$$\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'} \rightarrow \arg \min_{c'} \|\mathbf{h} - \boldsymbol{\mu}_{c'}\|_2$$

# Intermediate layers for generalization theory?

- How about intermediate layers? Do we have neural collapse?
- Empirical work by [GGB20] demonstrates existence of *effective depth*, which is a threshold  $L$  ---below layer  $L$  within-class variances decrease but no collapse, above layer  $L$  there is neural collapse
- Can we decompose overparametrized deep NNs into “representation learning component” and “interpolation component”? If so, helpful for generalization & transfer learning



Thank you!

**Contact:** [yiqiao.zhong@wisc.edu](mailto:yiqiao.zhong@wisc.edu)

**Office:** MSC 1122