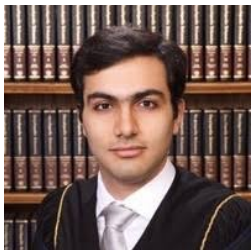


# Overparametrization in machine learning: insights from linear models

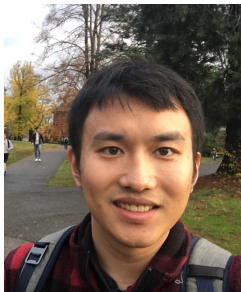
Andrea Montanari

Stanford University

March 16, 2023



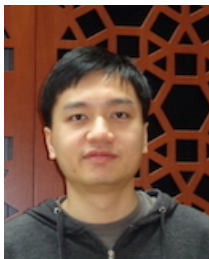
Behrooz Ghorbani



Song Mei



Theodor Misiakiewicz

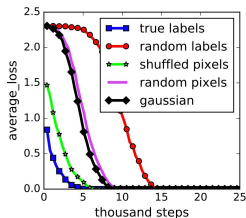


Yiqiao Zhong

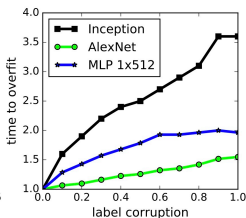


Chen Cheng

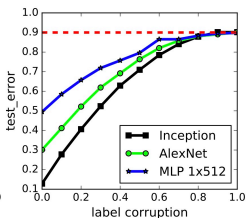
# Surprise #1: Near interpolation in ML practice



(a) learning curves



(b) convergence slowdown

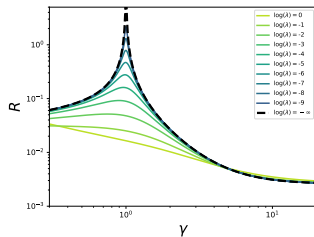
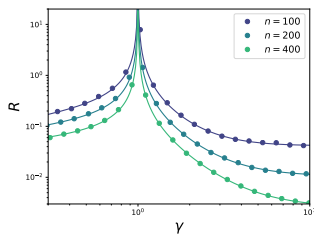


(c) generalization error growth

- ▶ Model complex enough to ‘interpolate’ random labels
- ▶ Despite this, does well on uncorrupted test samples
- ▶ Test error  $\gg$  Train error  $\approx 0$

## Surprise #2: Completely general

Test error of ridge(less) regression vs  $\gamma = p/n$



$$y_i = \langle \theta, x_i \rangle + \varepsilon_i, \quad x_i \sim N(0, I_d),$$
$$z_i = W^T x_i + g_i, \quad W \in \mathbb{R}^{d \times p}, \quad g_i \sim N(0, I_d).$$

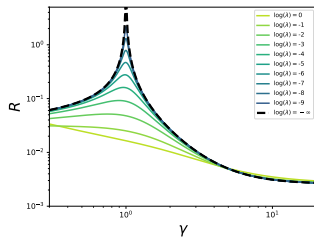
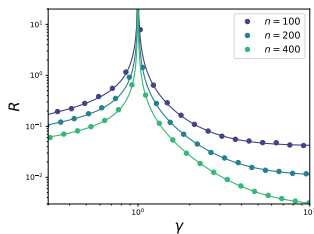
►  $x$ : latent.

$z$ : features

Regress  $y$  vs  $z$

## Surprise #2: Completely general

Test error of ridge(less) regression vs  $\gamma = p/n$



$$y_i = \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \varepsilon_i, \quad \mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d),$$
$$\mathbf{z}_i = \mathbf{W}^T \mathbf{x}_i + \mathbf{g}_i, \quad \mathbf{W} \in \mathbb{R}^{d \times p}, \quad \mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d).$$

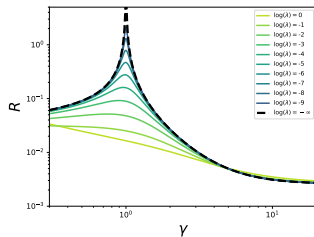
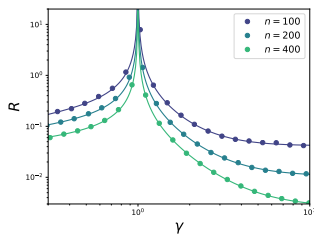
►  $\mathbf{x}$ : latent.

$\mathbf{z}$ : features

Regress  $y$  vs  $\mathbf{z}$

## Surprise #2: Completely general

Test error of ridge(less) regression vs  $\gamma = p/n$



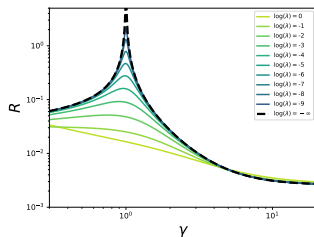
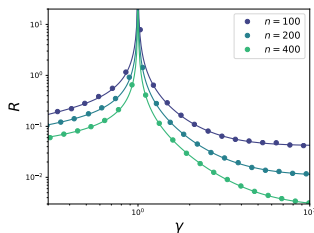
$$\begin{aligned} y_i &= \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \varepsilon_i, & \mathbf{x}_i &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \\ \mathbf{z}_i &= \mathbf{W}^\top \mathbf{x}_i + \mathbf{g}_i, & \mathbf{W} &\in \mathbb{R}^{d \times p}, \quad \mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \end{aligned}$$

►  $\mathbf{x}$ : latent.

$\mathbf{z}$ : features

Regress  $y$  vs  $\mathbf{z}$

## Surprise #2: Completely general



$$\begin{aligned} \mathbf{y}_i &= \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \varepsilon_i, & \mathbf{x}_i &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \\ \mathbf{z}_i &= \mathbf{W}^\top \mathbf{x}_i + \mathbf{g}_i, & \mathbf{W} &\in \mathbb{R}^{d \times p}, \quad \mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \end{aligned}$$

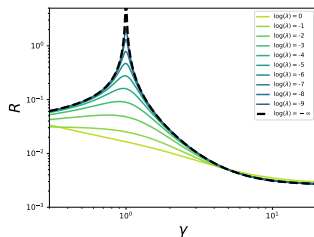
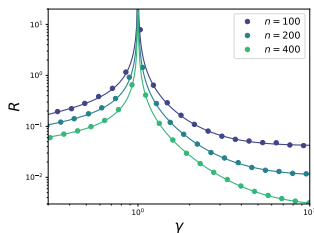
### Equivalent description

$$\begin{aligned} \mathbf{y}_i &= \langle \boldsymbol{\beta}, \mathbf{z}_i \rangle + \tilde{\varepsilon}_i, & \mathbf{z}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_d), \\ \boldsymbol{\Sigma} &= \mathbf{W}\mathbf{W}^\top + \mathbf{I}_p, & \boldsymbol{\beta} &\in \text{span}(\mathbf{W}). \end{aligned}$$

General picture?      Connection to neural nets?



## Surprise #2: Completely general



$$\begin{aligned} \mathbf{y}_i &= \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \varepsilon_i, & \mathbf{x}_i &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \\ \mathbf{z}_i &= \mathbf{W}^\top \mathbf{x}_i + \mathbf{g}_i, & \mathbf{W} &\in \mathbb{R}^{d \times p}, \quad \mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \end{aligned}$$

### Equivalent description

$$\begin{aligned} \mathbf{y}_i &= \langle \boldsymbol{\beta}, \mathbf{z}_i \rangle + \tilde{\varepsilon}_i, & \mathbf{z}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_d), \\ \boldsymbol{\Sigma} &= \mathbf{W}\mathbf{W}^\top + \mathbf{I}_p, & \boldsymbol{\beta} &\in \text{span}(\mathbf{W}). \end{aligned}$$

General picture?      Connection to neural nets?

# Outline

- 1 Examples of linear regression
- 2 A general formula
- 3 Benign overfitting
- 4 Kernel ridge regression
- 5 Random features
- 6 Neural tangent
- 7 Conclusion

## Examples of linear regression

# Setting

## ► Data

$$(\mathbf{y}_1, \mathbf{z}_1), (\mathbf{y}_2, \mathbf{z}_2), \dots, (\mathbf{y}_n, \mathbf{z}_n)$$

$\mathbf{y}_i \in \mathbb{R}$  : 'label', 'response',

$\mathbf{z}_i \in \mathbb{R}^p$  : 'features' vector', 'covariates'.

## ► Distribution

$$y_i = \langle \boldsymbol{\beta}, \mathbf{z}_i \rangle + \varepsilon_i, \quad \mathbb{E}(\varepsilon_i) = 0, \quad \mathbb{E}(\varepsilon_i^2) = \tau^2$$

# Setting

► Data

$$(\mathbf{y}_1, \mathbf{z}_1), (\mathbf{y}_2, \mathbf{z}_2), \dots, (\mathbf{y}_n, \mathbf{z}_n)$$

$\mathbf{y}_i \in \mathbb{R}$  : 'label', 'response',

$\mathbf{z}_i \in \mathbb{R}^p$  : 'features' vector', 'covariates'.

► Distribution

$$\mathbf{y}_i = \langle \boldsymbol{\beta}, \mathbf{z}_i \rangle + \varepsilon_i, \quad \mathbb{E}(\varepsilon_i) = 0, \quad \mathbb{E}(\varepsilon_i^2) = \tau^2$$

## Ridge regression

$$\hat{\boldsymbol{\beta}}(\lambda) := \arg \min_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}, \quad \mathbf{Z} := \begin{bmatrix} -\mathbf{z}_1- \\ \vdots \\ -\mathbf{z}_n- \end{bmatrix} \in \mathbb{R}^{p \times p}$$

$$\hat{\boldsymbol{\beta}}(\lambda) := \frac{1}{n} \mathbf{Z}^T (\mathbf{K}_n + (\lambda/n) \mathbf{I}_n)^{-1} \mathbf{y},$$
$$\mathbf{K}_n := \frac{1}{n} \mathbf{Z} \mathbf{Z}^T.$$

## Ridge regression

$$\hat{\boldsymbol{\beta}}(\lambda) := \arg \min_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}, \quad \mathbf{Z} := \begin{bmatrix} -\mathbf{z}_1- \\ \vdots \\ -\mathbf{z}_n- \end{bmatrix} \in \mathbb{R}^{p \times p}$$

$$\hat{\boldsymbol{\beta}}(\lambda) := \frac{1}{n} \mathbf{Z}^T (\mathbf{K}_n + (\lambda/n) \mathbf{I}_n)^{-1} \mathbf{y},$$
$$\mathbf{K}_n := \frac{1}{n} \mathbf{Z} \mathbf{Z}^T.$$

## Test error

$$\hat{\boldsymbol{\beta}}(\lambda) := \arg \min_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}.$$

$$\mathcal{R}_{\mathbf{Z}}(\lambda) := \mathbb{E}_{\text{new}} \left\{ (\mathbf{y}_{\text{new}} - \langle \hat{\boldsymbol{\beta}}(\lambda), \mathbf{z}_{\text{new}} \rangle)^2 \right\} - \underbrace{\mathbb{E}_{\text{new}} \left\{ (\mathbf{y}_{\text{new}} - \langle \boldsymbol{\beta}, \mathbf{z}_{\text{new}} \rangle)^2 \right\}}_{\text{Bayes}}$$

$$= \|\hat{\boldsymbol{\beta}}(\lambda) - \boldsymbol{\beta}\|_{\Sigma}^2 \quad \Sigma := \mathbb{E}[\mathbf{z}\mathbf{z}^T]$$



## Test error

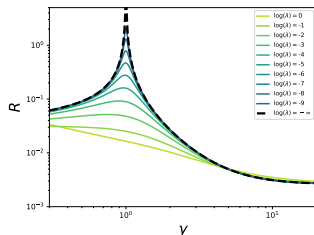
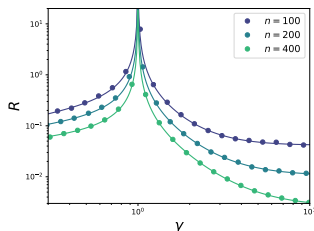
$$\hat{\boldsymbol{\beta}}(\lambda) := \arg \min_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}.$$

$$\mathcal{R}_{\mathbf{Z}}(\lambda) := \mathbb{E}_{\text{new}} \left\{ (\mathbf{y}_{\text{new}} - \langle \hat{\boldsymbol{\beta}}(\lambda), \mathbf{z}_{\text{new}} \rangle)^2 \right\} - \underbrace{\mathbb{E}_{\text{new}} \left\{ (\mathbf{y}_{\text{new}} - \langle \boldsymbol{\beta}, \mathbf{z}_{\text{new}} \rangle)^2 \right\}}_{\text{Bayes}}$$

$$= \|\hat{\boldsymbol{\beta}}(\lambda) - \boldsymbol{\beta}\|_{\boldsymbol{\Sigma}}^2 \quad \boldsymbol{\Sigma} := \mathbb{E}[\mathbf{z}\mathbf{z}^T]$$

## Examples

# Example #1: Well-concentrated covariates



▶  $\mathbf{z}_i = \Sigma^{1/2} \mathbf{x}_i$ ,  $\mathbb{E}\{\mathbf{x}_i \mathbf{x}_i^T\} = \mathbf{I}_p$ .

▶ Concentration properties for  $\mathbf{x}_i$ :

Either: Independent sub-Gaussian coordinates

or: Log-Sobolev

or: ...

## Example #2: Kernel Ridge Regression

### Data

$$\begin{aligned}\mathbf{x}_i &\sim \mathbb{P} \in \mathcal{P}(\mathbb{R}^d), \\ \mathbf{y}_i &= f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{R}^d; \mathbb{P}),\end{aligned}$$

### Function space view

$$\hat{f}_\lambda = \operatorname{argmin}_f \left\{ \sum_{i=1}^n (\mathbf{y}_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\},$$

$\mathcal{H} =$  Reproducing Kernel Hilbert Space.      Kernel  $\mathbf{K}$

## Example #2: Kernel Ridge Regression (Take 2)

### Data

$$\mathbf{x}_i \sim \mathbb{P} \in \mathcal{P}(\mathbb{R}^d),$$

$$\mathbf{y}_i = \mathbf{f}_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbf{f}_* \in \mathcal{H} \subseteq L^2(\mathbb{P}), \quad (\text{Hilbert space})$$

### Featurization map

$$\Phi : \mathbb{R}^d \rightarrow \mathcal{H}_0, \quad \mathbf{x} \mapsto \Phi(\mathbf{x}).$$

### Feature space view ( $p = \infty$ )

$$\hat{f}_\lambda(\mathbf{x}) = \langle \hat{\beta}_\lambda, \Phi(\mathbf{x}) \rangle, \quad \mathbf{z}_i = \Phi(\mathbf{x}_i)$$

$$\hat{\beta}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_{\mathcal{H}_0}^2 \right\},$$

---

$$K(\mathbf{x}_1, \mathbf{x}_2) = \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}_0}.$$

## Example #2: Kernel Ridge Regression (Take 2)

### Data

$$\mathbf{x}_i \sim \mathbb{P} \in \mathcal{P}(\mathbb{R}^d),$$

$$\mathbf{y}_i = \mathbf{f}_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbf{f}_* \in \mathcal{H} \subseteq L^2(\mathbb{P}), \quad (\text{Hilbert space})$$

### Featurization map

$$\boldsymbol{\Phi} : \mathbb{R}^d \rightarrow \mathcal{H}_0, \quad \mathbf{x} \mapsto \boldsymbol{\Phi}(\mathbf{x}).$$

### Feature space view ( $p = \infty$ )

$$\hat{\mathbf{f}}_\lambda(\mathbf{x}) = \langle \hat{\boldsymbol{\beta}}_\lambda, \boldsymbol{\Phi}(\mathbf{x}) \rangle, \quad \mathbf{z}_i = \boldsymbol{\Phi}(\mathbf{x}_i)$$

$$\hat{\boldsymbol{\beta}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_{\mathcal{H}_0}^2 \right\},$$

---

$$\mathbf{K}(\mathbf{x}_1, \mathbf{x}_2) = \langle \boldsymbol{\Phi}(\mathbf{x}_1), \boldsymbol{\Phi}(\mathbf{x}_2) \rangle_{\mathcal{H}_0}.$$

## Example #3: Random Features Regression

### Data

$$\begin{aligned}\mathbf{x}_i &\sim \mathbb{P} \in \mathcal{P}(\mathbb{R}^d), \\ y_i &= f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{P}),\end{aligned}$$

### Two-layer network with random first layer ( $p = N$ )

$$\hat{f}(\mathbf{x}; \mathbf{b}) = \sum_{i=1}^N b_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \mathbf{w}_i \sim \text{Unif}(\mathbb{S}^{d-1}),$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \mathbf{b}))^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$

## Example #3: Random Features Regression (Take 2)

### Data

$$\begin{aligned}\mathbf{x}_i &\sim \mathbb{P} \in \mathcal{P}(\mathbb{R}^d), \\ y_i &= f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{P}),\end{aligned}$$

Two-layer network with random first layer ( $p = N$ )

$$\hat{f}(\mathbf{x}; \mathbf{b}) = \sum_{i=1}^N b_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \mathbf{w}_i \sim \text{Unif}(\mathbb{S}^{d-1}),$$

$$\mathbf{z}_i = \Phi_{\mathbf{W}}(\mathbf{x}_i) := (\sigma(\langle \mathbf{w}_1, \mathbf{x}_i \rangle), \sigma(\langle \mathbf{w}_2, \mathbf{x}_i \rangle), \dots, \sigma(\langle \mathbf{w}_N, \mathbf{x}_i \rangle))^T,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$



## Example #4: (Neural) Tangent Regression

- ▶ Parametric model  $\alpha f(\cdot; \boldsymbol{\theta}) : \mathbb{R}^d \rightarrow \mathbb{R}$  (parameters:  $(\alpha, \boldsymbol{\theta})$ )
- ▶ Linearize around SGD initialization  $\boldsymbol{\theta}^0$

$$\begin{aligned}\alpha f(\mathbf{x}; \boldsymbol{\theta}^0 + \alpha^{-1} \mathbf{b}) &= \alpha f(\mathbf{x}; \boldsymbol{\theta}^0) + \langle \mathbf{b}, \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}^0) \rangle + O(\alpha^{-1}) \\ &= \text{const.} + \underbrace{\langle \mathbf{b}, \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}^0) \rangle}_{f_{\text{NT}}(\mathbf{x}; \mathbf{b})} + O(\alpha^{-1})\end{aligned}$$

---

Jacot, Gabriel, Hongler, 2018; Du, Zhai, Póczos, Singh 2018; Allen-Zhu, Li, Song 2018; Chizat, Bach, 2019; Ghorbani, Mei, Misiakiewicz, M, 2019; Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019; Oymak, Soltanolkotabi, 2019; . . .

## Example #4: (Neural) Tangent Regression

Two-layer neural net

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$f_{\text{NT}}(\mathbf{x}; \bar{\mathbf{b}}, \mathbf{b}) = \sum_{j=1}^N \langle \mathbf{b}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j^0, \mathbf{x} \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle \boldsymbol{\theta}_j^0, \mathbf{x} \rangle).$$

## Example #4: (Neural) Tangent Regression

### Two-layer neural net

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$f_{\text{NT}}(\mathbf{x}; \bar{\mathbf{b}}, \mathbf{b}) = \sum_{j=1}^N \langle \mathbf{b}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j^0, \mathbf{x} \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle \mathbf{w}_j^0, \mathbf{x} \rangle).$$

Two-layer network with random first layer ( $p = Nd$ )

$$\mathbf{z}_i = \Phi_{\mathbf{W}}(\mathbf{x}_i) := (\mathbf{x}_i^T \sigma'(\langle \mathbf{w}_1^0, \mathbf{x}_i \rangle), \mathbf{x}_i^T \sigma'(\langle \mathbf{w}_2^0, \mathbf{x}_i \rangle), \dots, \mathbf{x}_i^T \sigma'(\langle \mathbf{w}_N^0, \mathbf{x}_i \rangle))^T,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$

## Example #4: (Neural) Tangent Regression

Two-layer neural net

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$f_{\text{NT}}(\mathbf{x}; \bar{\mathbf{b}}, \mathbf{b}) = \sum_{j=1}^N \langle \mathbf{b}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j^0, \mathbf{x} \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle \mathbf{w}_j^0, \mathbf{x} \rangle).$$

Two-layer network with random first layer ( $p = Nd$ )

$$\mathbf{z}_i = \Phi_{\mathbf{W}}(\mathbf{x}_i) := (\mathbf{x}_i^T \sigma'(\langle \mathbf{w}_1^0, \mathbf{x}_i \rangle), \mathbf{x}_i^T \sigma'(\langle \mathbf{w}_2^0, \mathbf{x}_i \rangle), \dots, \mathbf{x}_i^T \sigma'(\langle \mathbf{w}_N^0, \mathbf{x}_i \rangle))^T,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$

A general formula

Assumptions:  $p = \infty$  ( $\boldsymbol{\beta}, \mathbf{z}_i \in \text{Hilbert}$ )

$$\mathbf{y}_i = \langle \boldsymbol{\beta}, \mathbf{z}_i \rangle + \varepsilon_i, \quad \mathbb{E}(\varepsilon_i) = \mathbb{E}(\varepsilon_i \mathbf{z}_i) = 0, \quad \mathbb{E}(\varepsilon_i^2) = \tau^2$$

1.  $\text{Tr}(\boldsymbol{\Sigma}) < \infty$  and (wlog)  $\|\boldsymbol{\Sigma}\| = 1$ .
2.  $\|\boldsymbol{\Sigma}^{-1/2} \boldsymbol{\beta}\| < \infty$ .
3.  $(\sigma_i)_{i \geq 1}$ : ordered eigenvalues of  $\boldsymbol{\Sigma}$ . For all  $1 \leq k \leq n$  :

$$\sum_{l=k}^{\infty} \sigma_l \leq d_{\boldsymbol{\Sigma}} \sigma_k.$$

4. For  $\mathbf{u}_i := \boldsymbol{\Sigma}^{-1/2} \mathbf{z}_i$  one of the following holds:
  - (a) Independent sub-Gaussian coordinates.
  - (b) Concentration 1-Lipschitz convex function (eg implied by Log-Sobolev)

## General result

Theorem (Cheng, M, 2022)

- ▶  $\mathcal{R}_{\mathbf{Z}}(\lambda)$  be the test error of ridge regression (conditional on  $\mathbf{Z}$ )
- ▶  $R_n^s(\lambda)$  (non-random) test error in the equivalent sequence model.

Then

$$\mathcal{R}_{\mathbf{Z}}(\lambda) = (1 + \text{err}_n)R_n^s(\lambda),$$

where  $\text{err}_n$  is small with high probability, provided ...

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’ [Most need a separate proof!]

## General result

Theorem (Cheng, M, 2022)

- ▶  $\mathcal{R}_{\mathbf{Z}}(\lambda)$  be the test error of ridge regression (conditional on  $\mathbf{Z}$ )
- ▶  $R_n^s(\lambda)$  (non-random) test error in the equivalent sequence model.

Then

$$\mathcal{R}_{\mathbf{Z}}(\lambda) = (1 + \text{err}_n)R_n^s(\lambda),$$

where  $\text{err}_n$  is small with high probability, provided ...

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’ [Most need a separate proof!]



# General result

## Theorem (Cheng, M, 2022)

- ▶  $\mathcal{R}_{\mathbf{Z}}(\lambda)$  be the test error of ridge regression (conditional on  $\mathbf{Z}$ )
- ▶  $R_n^s(\lambda)$  (non-random) test error in the equivalent sequence model.

Then

$$\mathcal{R}_{\mathbf{Z}}(\lambda) = (1 + \text{err}_n)R_n^s(\lambda),$$

where  $\text{err}_n$  is small with high probability, provided ...

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’ [Most need a separate proof!]

One key quantity controlling  $\text{err}_n$  ( $\lambda = 0+$ )

$$\chi_n := \frac{\sigma_{[\eta n]} d_{\Sigma} \log^2(d_{\Sigma})}{\kappa n \lambda_{\star}(0)},$$

▶  $\lambda_{\star} \asymp \sigma_{[cn]}$ .

▶ Need  $d_{\Sigma} \leq n^{1+\varepsilon}$

# The actual theorem

1. The ratio between effective dimension and regularization parameter:

$$\chi_n(\lambda) := 1 + \frac{\sigma_{|\eta n|} \mathbf{d}_\Sigma \log^2(\mathbf{d}_\Sigma)}{\lambda}. \quad (20)$$

Here  $\eta$  is a constant that only depends on  $\mathbf{C}_x$ , and hence we will leave it implicit.

2. The ratio between regularization and effective regularization

$$\kappa := \min\left(\frac{\lambda}{n\lambda_*}; 1 - \frac{\lambda}{n\lambda_*}\right) > 0. \quad (21)$$

3. For a positive semi-definite operator  $\mathbf{Q}$ , define the modified population resolvent:

$$\mathcal{R}_0(\mu_0, \mu; \mathbf{Q}) := \text{Tr}\left(\Sigma^{\frac{1}{2}} \mathbf{Q} \Sigma^{\frac{1}{2}} (\mu_0 \mathbf{I} + \mu \Sigma)^{-1}\right). \quad (22)$$

Letting  $\beta = \Sigma^{1/2} \boldsymbol{\theta}$ ,  $\|\boldsymbol{\theta}\| < \infty$ , we consider the ratio

$$\rho(\lambda) := \frac{\mathcal{R}_0(\lambda_*, 1; \boldsymbol{\theta} \boldsymbol{\theta}^\top / \|\boldsymbol{\theta}\|^2)}{\mathcal{R}_0(\lambda_*, 1; \mathbf{I})} \in (0, 1]. \quad (23)$$

We next present our master theorem for ridge regression: its proof is postponed to Section 6.

**Theorem 1** (Ridge regression). *Under Assumption 1, for any positive integers  $k$  and  $D$ , there exist constants  $\eta = \eta(\mathbf{C}_x) \in (0, 1/2)$  and  $\mathbf{C} = \mathbf{C}(\mathbf{C}_x, D) > 0$  such that the following hold. Define  $\chi_n(\lambda)$ ,  $\kappa$ ,  $\rho(\lambda)$  as above (with  $\eta = \eta(\mathbf{C}_x)$  in Eq. (20)).*

*If it holds that*

$$\chi_n(\lambda)^3 \log^2 n \leq \mathbf{C} n \kappa^{4.5}, \quad n^{-2D+1} = \mathcal{O}\left(\sqrt{\frac{\kappa^3 \log^2 n}{\max\{n, \lambda\}}}\right),$$

*then for all  $n = \Omega_{k,D}(1)$ , with probability  $1 - \mathcal{O}_k(n^{-D+1})$  we have:*

1. **Variance approximation.**

$$|\mathcal{Y}_X(\lambda) - \mathbf{V}_n(\lambda)| = \mathcal{O}_{k, \mathbf{C}_x, D}\left(\frac{\chi_n(\lambda)^3 \log^2 n}{n^{1-\frac{1}{k}} \kappa^{9.5}}\right) \cdot \mathbf{V}_n(\lambda).$$

2. **Bias approximation.** *If we additionally have  $\chi_n(\lambda)^3 \log^2 n \leq \mathbf{C} n \kappa^{4.5} \sqrt{\rho(\lambda)}$  and  $\lambda k n^{-\frac{1}{k}} \leq n \kappa / 2$ , for all  $n = \Omega_{k,D}(1)$ , we have*

$$|\mathcal{B}_X(\lambda) - \mathbf{B}_n(\lambda)| = \mathcal{O}_{k, \mathbf{C}_x, D}\left(\frac{\lambda_*(\lambda)^{k+1}}{n \kappa^3} + \frac{\chi_n(\lambda)^3 \log^2 n}{\sqrt{\rho(\lambda)} n^{1-\frac{1}{k}} \kappa^{8.5}}\right) \cdot \mathbf{B}_n(\lambda).$$

**Remark 3.1** The condition  $\|\boldsymbol{\theta}\| < \infty$  in Assumption 1 amounts to requiring that the coef

## Equivalent sequence model

$$\theta_i := \langle \boldsymbol{\beta}, \mathbf{v}_i \rangle, \quad \mathbf{v}_i := i\text{-th eigenvectors of } \boldsymbol{\Sigma}$$

$$\mathbf{y}_i^s = \sigma_i^{1/2} \theta_i + \frac{\omega}{\sqrt{n}} \mathbf{g}_i, \quad (\mathbf{g}_i)_{i \geq 1} \sim_{\text{iid}} \mathbf{N}(0, 1),$$

$$\hat{\theta}_i^s := \operatorname{argmin}_{t \in \mathbb{R}} \{ (\mathbf{y}_i^s - \sigma_i^{1/2} t)^2 + \lambda_\star t^2 \} = \frac{\sigma_i^{1/2}}{\sigma_i + \lambda_\star} \cdot \mathbf{y}_i^s.$$

Effective noise level and regularization  $\omega, \lambda_\star$

$$\omega^2 = \tau^2 + \underbrace{\mathbb{E}_{\mathbf{g}} \left\{ \sum_{i \geq 1} \sigma_i (\hat{\theta}_i^s - \theta_i)^2 \right\}}_{R_n^s}, \quad n - \frac{\lambda}{\lambda_\star} = \sum_{i \geq 1} \frac{\sigma_i}{\sigma_i + \lambda_\star}$$

## Equivalent sequence model

$$\theta_i := \langle \boldsymbol{\beta}, \mathbf{v}_i \rangle, \quad \mathbf{v}_i := i\text{-th eigenvectors of } \boldsymbol{\Sigma}$$

$$\mathbf{y}_i^s = \sigma_i^{1/2} \theta_i + \frac{\omega}{\sqrt{n}} \mathbf{g}_i, \quad (\mathbf{g}_i)_{i \geq 1} \sim_{\text{iid}} \mathbf{N}(0, 1),$$

$$\hat{\theta}_i^s := \operatorname{argmin}_{t \in \mathbb{R}} \{ (\mathbf{y}_i^s - \sigma_i^{1/2} t)^2 + \lambda_\star t^2 \} = \frac{\sigma_i^{1/2}}{\sigma_i + \lambda_\star} \cdot \mathbf{y}_i^s.$$

Effective noise level and regularization  $\omega, \lambda_\star$

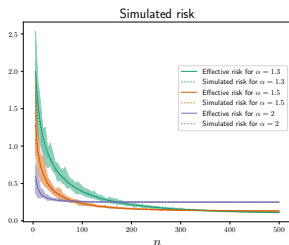
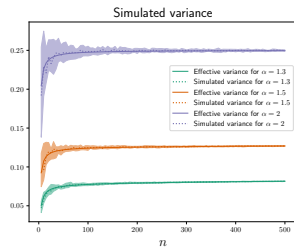
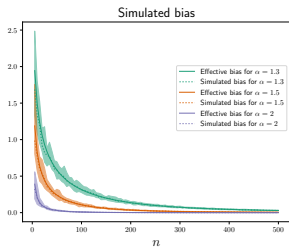
$$\omega^2 = \tau^2 + \underbrace{\mathbb{E}_{\mathbf{g}} \left\{ \sum_{i \geq 1} \sigma_i (\hat{\theta}_i^s - \theta_i)^2 \right\}}_{R_n^s}, \quad n - \frac{\lambda}{\lambda_\star} = \sum_{i \geq 1} \frac{\sigma_i}{\sigma_i + \lambda_\star}$$

## Specific eigenvalue structures

**Power law decay:**  $\sigma_i = i^{-\alpha}, \alpha > 1$

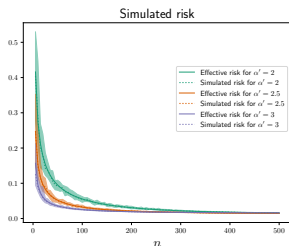
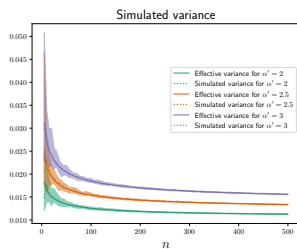
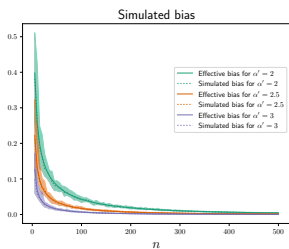
**Critical decay:**  $\sigma_i = i^{-1}(1 + \log i)^{-\alpha'}$ .

# Power law decay; $\lambda = 0+$



- ▶ Variance does not decrease with  $n$ .
- ▶ Need to use larger  $\lambda$ .

# Critical decay; $\lambda = 0+$



- ▶ Variance does not decrease with  $n$ !
- ▶ Benign overfitting

[Bartlett, Long, Lugosi, Tsigler, 2020]



## Benign overfitting

## Simplifying formulas in the seq. model

Determine

$$n - \frac{\lambda}{\lambda_{\star}} = \text{Tr} \left( \boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-1} \right)$$

Then  $R_n^s(\lambda) = B_n^s(\lambda) + V_n^s(\lambda)$ :

$$B_n^s(\lambda) = \frac{\lambda_{\star}^2 \langle \boldsymbol{\beta}, (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-2} \boldsymbol{\Sigma} \boldsymbol{\beta} \rangle}{1 - n^{-1} \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-2} \right)},$$
$$V_n^s(\lambda) = \frac{\tau^2 \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-2} \right)}{n - \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_{\star} \mathbf{I})^{-2} \right)}.$$

## Eigenvalue decay assumption

$$\mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq \mathrm{Tr} \left( \boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1} \right) = n - \frac{\lambda}{\lambda_*}$$

Assume:  $\mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq n(1 - c_*^{-1})$

Then

$$\begin{aligned} V_n^s(\lambda) &= \frac{\tau^2 \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)} \\ &\leq \frac{c_* \tau^2}{n} \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \end{aligned}$$

## Eigenvalue decay assumption

$$\text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq \text{Tr} \left( \boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1} \right) = n - \frac{\lambda}{\lambda_*}$$

Assume:  $\text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq n(1 - c_*^{-1})$

Then

$$\begin{aligned} V_n^s(\lambda) &= \frac{\tau^2 \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)}{n - \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)} \\ &\leq \frac{c_* \tau^2}{n} \text{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \end{aligned}$$

## Eigenvalue decay assumption

$$\mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq \mathrm{Tr} \left( \boldsymbol{\Sigma} (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-1} \right) = n - \frac{\lambda}{\lambda_*}$$

Assume:  $\mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \leq n(1 - c_*^{-1})$

Then

$$\begin{aligned} V_n^s(\lambda) &= \frac{\tau^2 \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)}{n - \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right)} \\ &\leq \frac{c_* \tau^2}{n} \mathrm{Tr} \left( \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \end{aligned}$$

Continuing...

$$\begin{aligned} V_n^s(\lambda) &\leq \frac{c_*\tau^2}{n} \text{Tr} \left( \mathbf{\Sigma}^2 (\mathbf{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \\ &\leq \frac{c_*\tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_*\tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_*\tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\sigma_{k_*+1}^2} \right\} \end{aligned}$$

For  $k_* := \max\{k : k \geq \lambda_*\}$

Continuing...

$$\begin{aligned} V_n^s(\lambda) &\leq \frac{c_* \tau^2}{n} \text{Tr} \left( \mathbf{\Sigma}^2 (\mathbf{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\sigma_{k_*+1}^2} \right\} \end{aligned}$$

For  $k_* := \max\{k : k \geq \lambda_*\}$

Continuing...

$$\begin{aligned} V_n^s(\lambda) &\leq \frac{c_* \tau^2}{n} \text{Tr} \left( \mathbf{\Sigma}^2 (\mathbf{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\sigma_{k_*+1}^2} \right\} \end{aligned}$$

For  $k_* := \max\{k : k \geq \lambda_*\}$



Continuing...

$$\begin{aligned} V_n^s(\lambda) &\leq \frac{c_* \tau^2}{n} \text{Tr} \left( \mathbf{\Sigma}^2 (\mathbf{\Sigma} + \lambda_* \mathbf{I})^{-2} \right) \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\lambda_*^2} \right\} \\ &\leq \frac{c_* \tau^2}{n} \left\{ k_* + \sum_{\ell=k_*+1}^{\infty} \frac{\sigma_\ell^2}{\sigma_{k_*+1}^2} \right\} \end{aligned}$$

For  $k_* := \max\{k : k \geq \lambda_*\}$

# Benign overfitting

## Proposition (Cheng, M, 2022)

Let  $k_\star := \max\{k : \sigma_k \geq \lambda_\star\}$ ,  $b_k := \sigma_k/\sigma_{k+1}$  and

$$r_q(k) := \sum_{\ell > k} \left( \frac{\sigma_\ell}{\sigma_{k+1}} \right)^q, \quad \bar{r}(k) := \frac{r_1(k)^2}{r_2(k)}.$$

Then,

$$V_n(\lambda) \leq c_\star \tau^2 \left( \frac{k_\star}{n} + \frac{r_2(k_\star)}{n} \right) \leq c_\star \tau^2 \left( \frac{k_\star}{n} + \frac{4b_{k_\star}^2 n}{\bar{r}(k_\star)} \right),$$
$$B_n(\lambda) \leq c_\star \left( \sigma_{k_\star}^2 \|\boldsymbol{\beta}_{\leq k_\star}\|_{\Sigma^{-1}}^2 + \|\boldsymbol{\beta}_{> k_\star}\|_{\Sigma}^2 \right).$$

- ▶ Consistent if:
  - ▶  $1 \ll k_\star \ll n$ .
  - ▶  $\bar{r}(k_\star) \rightarrow \infty$
  - ▶  $\|\boldsymbol{\beta}_{> k_\star}\|_{\Sigma}^2 \rightarrow 0$
- ▶ cf. Bartlett, Long, Lugosi, Tsigler, 2020; Bartlett, Tsigler, 2021

## Kernel ridge regression

## High-dimensional setting

- ▶  $\mathbf{x}_i \sim \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$
- ▶  $\mathbf{y}_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in L^2(\mathbb{R}^d; \mathbb{P})$
- ▶  $K(\mathbf{x}_1, \mathbf{x}_2) = h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d), \quad \mathbb{E}[h(\mathbf{G})\text{He}_k(\mathbf{G})] \neq 0$  for all  $k$ .

$$\hat{f}_\lambda = \operatorname{argmin}_f \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}.$$

# Staircase phenomenon

Theorem (Ghorbani, Mei, Misiakiewicz, M. 2019)

Let  $\ell \in \mathbb{Z}$ , and assume  $d^{\ell+\varepsilon} \leq n \leq d^{\ell+1-\varepsilon}$ ,  $\varepsilon > 0$ . Then, for any  $\lambda \in [0, \lambda_*(\sigma)]$ ,

$$R_\infty(f_*; \lambda) = \|P_{>\ell} f_*\|_{L^2}^2 + \|f_*\|_{L^2}^2 o_d(1),$$

$P_{>\ell} f_* =$  Projection of  $f_*$  onto deg.  $> \ell$  polynomials

*Further, no inner product kernel method can do better.*

- ▶ Pointwise result (valid any for fixed  $f_*$ )
- ▶  $\lambda = 0$ : optimal (interpolation)

---

Liang, Rakhlin, Zhai, 2019:  $\|f_*\|_{\mathcal{K}} \leq C$ , upper bounds on the rates. Bartlett, Rakhlin, M., 2021: Sharp results for  $n \asymp d$ .

# Staircase phenomenon

Theorem (Ghorbani, Mei, Misiakiewicz, M. 2019)

Let  $\ell \in \mathbb{Z}$ , and assume  $d^{\ell+\varepsilon} \leq n \leq d^{\ell+1-\varepsilon}$ ,  $\varepsilon > 0$ . Then, for any  $\lambda \in [0, \lambda_*(\sigma)]$ ,

$$R_\infty(f_*; \lambda) = \|P_{>\ell} f_*\|_{L^2}^2 + \|f_*\|_{L^2}^2 o_d(1),$$

$P_{>\ell} f_* = \text{Projection of } f_* \text{ onto deg. } > \ell \text{ polynomials}$

Further, no inner product kernel method can do better.

- ▶ Pointwise result (valid any for fixed  $f_*$ )
- ▶  $\lambda = 0$ : optimal (interpolation)

---

Liang, Rakhlin, Zhai, 2019:  $\|f_*\|_{\mathcal{K}} \leq C$ , upper bounds on the rates. Bartlett, Rakhlin, M., 2021: Sharp results for  $n \asymp d$ .

# Staircase phenomenon

Theorem (Ghorbani, Mei, Misiakiewicz, M. 2019)

Let  $\ell \in \mathbb{Z}$ , and assume  $d^{\ell+\varepsilon} \leq n \leq d^{\ell+1-\varepsilon}$ ,  $\varepsilon > 0$ . Then, for any  $\lambda \in [0, \lambda_*(\sigma)]$ ,

$$R_\infty(f_*; \lambda) = \|P_{>\ell} f_*\|_{L^2}^2 + \|f_*\|_{L^2}^2 o_d(1),$$

$P_{>\ell} f_* =$  Projection of  $f_*$  onto deg.  $> \ell$  polynomials

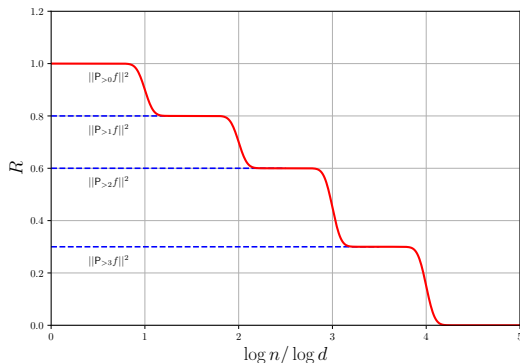
Further, no inner product kernel method can do better.

- ▶ Pointwise result (valid any for fixed  $f_*$ )
- ▶  $\lambda = 0$ : optimal (interpolation)

---

Liang, Rakhlin, Zhai, 2019:  $\|f_*\|_{\mathcal{K}} \leq C$ , upper bounds on the rates. Bartlett, Rakhlin, M., 2021: Sharp results for  $n \asymp d$ .

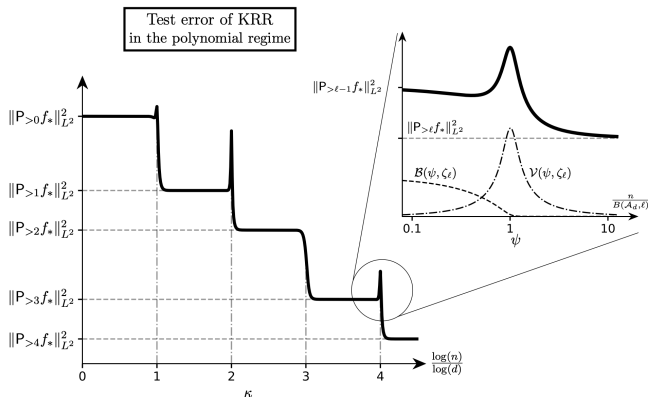
# Sketch



- ▶ If  $n \leq d^{1.99}$  can fit only linear functions.
- ▶ Valid for any inner product kernel
- ▶ Includes fully connected multi-layer nets



# Follow-up work



[Misiakiewicz, 4/2022; Xiao, Pennington, 5/2022; Hu, Lu, 5/2022 ]

## Intuition

$$K(\mathbf{x}_1, \mathbf{x}_2) = h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d):$$

$$K(\mathbf{x}_1, \mathbf{x}_2) = h_{\leq \ell}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) + h_{> \ell}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d)$$

$$\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}:$$

$$\mathbf{K}_n \approx \mathbf{Y}_{\leq \ell} \mathbf{D} \mathbf{Y}_{\leq \ell}^T + \lambda_* \mathbf{I}_n$$

## Intuition

$$K(\mathbf{x}_1, \mathbf{x}_2) = h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d):$$

$$K(\mathbf{x}_1, \mathbf{x}_2) = h_{\leq \ell}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) + h_{> \ell}(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d)$$

$$\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}:$$

$$\mathbf{K}_n \approx \underbrace{\mathbf{Y}_{\leq \ell} \mathbf{D} \mathbf{Y}_{\leq \ell}^T}_{\text{low rank, large norm}} + \underbrace{\lambda_* \mathbf{I}_n}_{\text{self-induced regularization}}$$

## Random features

## High-dimensional setting

- ▶  $\mathbf{x}_i \sim \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$
- ▶  $\mathbf{y}_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in L^2(\mathbb{R}^d; \mathbb{P})$

$$\hat{f}(\mathbf{x}; \mathbf{b}) = \sum_{i=1}^N b_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle), \quad \mathbf{w}_i \sim \text{Unif}(\mathbb{S}^{d-1}),$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \mathbf{b}))^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$

# Proportional regime

$$n, N, d \rightarrow \infty$$

$$\frac{N}{d} \rightarrow \psi_1, \quad \frac{n}{d} \rightarrow \psi_2.$$

# Precise asymptotics

## Theorem (Mei, M. 2019)

Decompose  $\sigma(\mathbf{x}) = \sigma_0 + \sigma_1 \mathbf{x} + \sigma^{\text{NL}}(\mathbf{x})$  where (for  $\mathbf{G} \sim \mathbf{N}(0, \mathbf{1})$ )

$$\mathbb{E}[\mathbf{G} \sigma^{\text{NL}}(\mathbf{G})] = \mathbb{E}[\sigma^{\text{NL}}(\mathbf{G})] = 0, \quad \zeta^2 := \frac{\sigma_1^2}{\mathbb{E}[\sigma^{\text{NL}}(\mathbf{G})^2]}.$$

Then, for any  $\bar{\lambda} = \lambda / \bar{b}_*^2 > 0$

$$R(\hat{f}_\lambda) = F_1^2 \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) + (\tau^2 + F_*^2) \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) + F_*^2 + o_d(1),$$

where  $\mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda})$ ,  $\mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda})$  are explicitly given below.

# Precise asymptotics

## Theorem (Mei, M. 2019)

Decompose  $\sigma(\mathbf{x}) = \sigma_0 + \sigma_1 \mathbf{x} + \sigma^{\text{NL}}(\mathbf{x})$  where (for  $\mathbf{G} \sim \mathbf{N}(0, \mathbf{1})$ )

$$\mathbb{E}[\mathbf{G} \sigma^{\text{NL}}(\mathbf{G})] = \mathbb{E}[\sigma^{\text{NL}}(\mathbf{G})] = 0, \quad \zeta^2 := \frac{\sigma_1^2}{\mathbb{E}[\sigma^{\text{NL}}(\mathbf{G})^2]}.$$

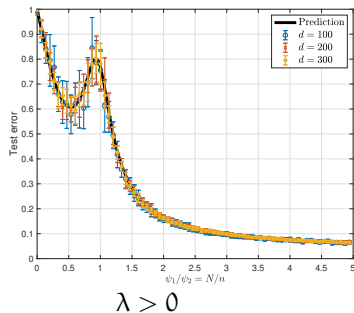
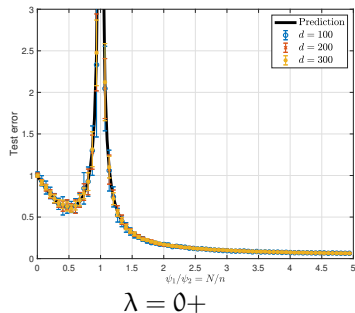
Then, for any  $\bar{\lambda} = \lambda / \bar{\mathbf{b}}_*^2 > 0$

$$\mathbf{R}(\hat{\mathbf{f}}_\lambda) = F_1^2 \mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda}) + (\tau^2 + F_*^2) \mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda}) + F_*^2 + \mathbf{o}_d(1),$$

where  $\mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda})$ ,  $\mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda})$  are explicitly given below.



# Risk vs overparametrization



- ▶ Solid line: Theoretical prediction (Random matrix theory)

## Neural tangent

## Neural Tangent: Parameters view

$$\hat{f}_{\text{NT}}(\mathbf{x}; \mathbf{b}) = \sum_{j=1}^N \langle \hat{\mathbf{b}}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j, \mathbf{x} \rangle).$$

$$\mathbf{z}_i := (\mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_1, \mathbf{x}_i \rangle), \mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_2, \mathbf{x}_i \rangle), \dots, \mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_N, \mathbf{x}_i \rangle))^\top,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}.$$

## Neural Tangent: Function view

$$\hat{f}_{\text{NT}} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathbb{K}_N}^2 \right\}.$$

$\|\cdot\|_{\mathbb{K}_N}$ : RKHS norm

$$\mathbb{K}_N(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{Nd} \sum_{i=1}^N \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{x}_1 \rangle) \sigma'(\langle \mathbf{w}_i, \mathbf{x}_2 \rangle)$$

## Finite vs infinite width: Kernel view

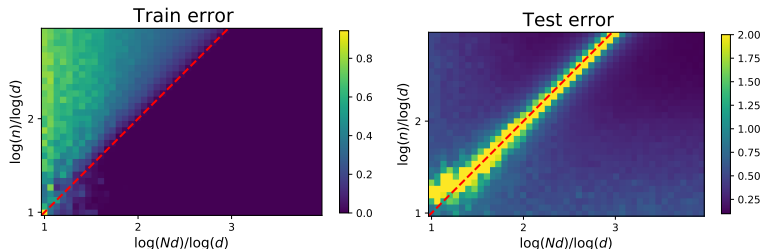
$$K_N(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{Nd} \sum_{i=1}^N \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{x}_1 \rangle) \sigma'(\langle \mathbf{w}_i, \mathbf{x}_2 \rangle)$$
$$K(\mathbf{x}_1, \mathbf{x}_2) = \frac{\langle \mathbf{x}_1, \mathbf{x}_2 \rangle}{d} E_{\mathbf{w}} \{ \sigma'(\langle \mathbf{w}, \mathbf{x}_1 \rangle) \sigma'(\langle \mathbf{w}, \mathbf{x}_2 \rangle) \}$$

$$K_N \rightarrow K$$

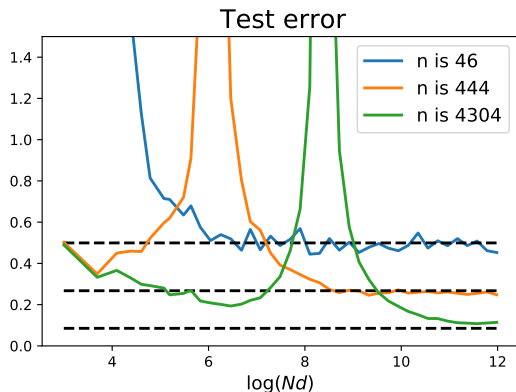
- ▶ What notion of convergence?
- ▶ How big should N be?

# Another small experiment

( $d = 20$ )



- ▶  $f_*(\mathbf{x}) = g(\langle \boldsymbol{\beta}_*, \mathbf{x} \rangle)$ ,  $g = \text{deg-4 polynomial}$



- ▶  $f_*(\mathbf{x}) = g(\langle \boldsymbol{\beta}_*, \mathbf{x} \rangle)$ ,  $g = \text{deg-4 polynomial}$
- ▶ NT ridge regression vs Kernel Ridge(-less) Regression

$$\hat{f} := \arg \min_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \left\{ \|f\|_K \text{ subj. to } f(\mathbf{x}_i) = y_i \forall i \leq n \right\},$$

## Rigorous confirmation

$\mathcal{R}_N(\mathbf{f}_*; \lambda) :=$  Risk of linearized network.

$\mathcal{R}_\infty(\mathbf{f}_*; \lambda) :=$  Risk of KRR.

Theorem (M, Zhong, 2020, 2021)

Assume  $d^\ell \ll n \ll d^{\ell+1}$  for some integer  $\ell$ . Then

$$\mathcal{R}_N(\mathbf{f}_*; \lambda) = \mathcal{R}_\infty(\mathbf{f}_*; \lambda) + O\left(\|\mathbf{f}_*\|_{L^2}^2 \sqrt{\frac{n(\log n)^c}{Nd}}\right)$$



# Insights

$$R_N(f_*; \lambda) = R_\infty(f_*; \lambda) + O\left(\sqrt{\frac{n(\log n)^C}{Nd}}\right)$$

**Insight 1:** Risk constant for  $Nd \gtrsim n(\log n)^C$

**Insight 2:** Overparametrization does not hurt

**Insight 3:** Interpolation ( $\lambda = 0$ ) nearly optimal (see KRR result)

Intuition for 3?

# Insights

$$R_N(f_*; \lambda) = R_\infty(f_*; \lambda) + O\left(\sqrt{\frac{n(\log n)^C}{Nd}}\right)$$

**Insight 1:** Risk constant for  $Nd \gtrsim n(\log n)^C$

**Insight 2:** Overparametrization does not hurt

**Insight 3:** Interpolation ( $\lambda = 0$ ) nearly optimal (see KRR result)

**Intuition for 3?**

## Intuition

- ▶ High-degree part of  $\sigma \approx$  Noise in the features
- ▶ Noise in the features  $\approx$  Diagonal term in the  $\mathbf{K}$
- ▶ Diagonal term in  $\mathbf{K} \approx$  Non-zero ridge regularization

## Conclusion

# Conclusion

- ▶ Linear models elucidate generalization puzzle in deep learning!
- ▶ Neural tangent model: optimal interpolation at high SNR
- ▶ Towards a unified theory

Thanks!

## Conclusion

- ▶ Linear models elucidate generalization puzzle in deep learning!
- ▶ Neural tangent model: optimal interpolation at high SNR
- ▶ Towards a unified theory

Thanks!

## Conclusion

- ▶ Linear models elucidate generalization puzzle in deep learning!
- ▶ Neural tangent model: optimal interpolation at high SNR
- ▶ Towards a unified theory

Thanks!

## Conclusion

- ▶ Linear models elucidate generalization puzzle in deep learning!
- ▶ Neural tangent model: optimal interpolation at high SNR
- ▶ Towards a unified theory

Thanks!