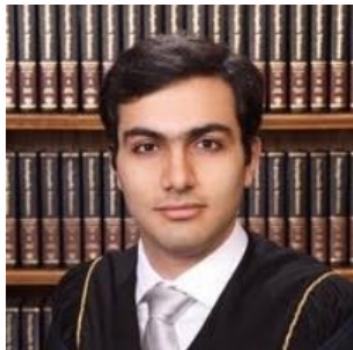


# Overparametrized models: Linear theory and its limits

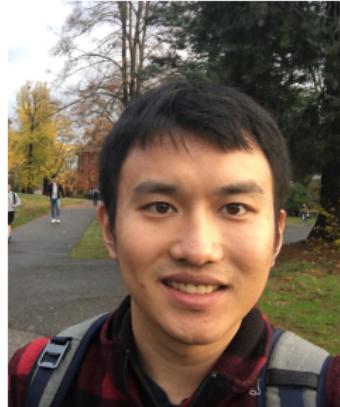
Andrea Montanari

Stanford University

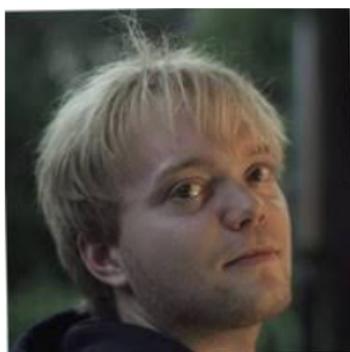
August 12, 2025



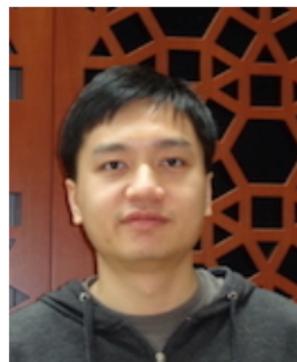
Behrooz Ghorbani



Song Mei



Theodor Misiakiewicz

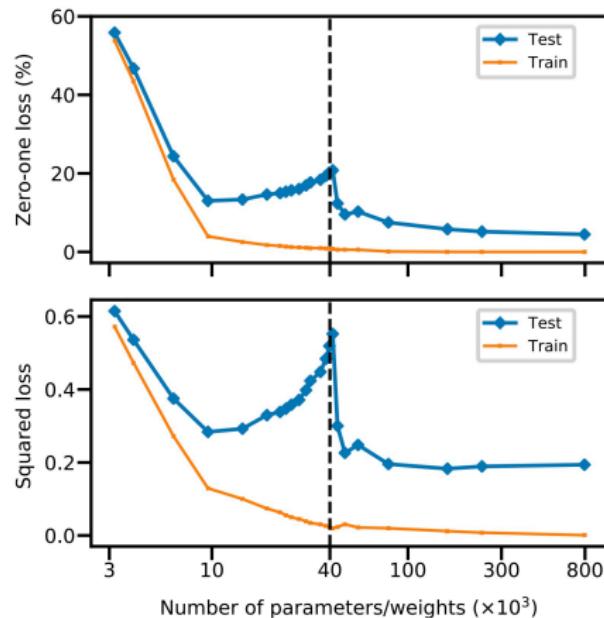


Yiqiao Zhong



Chen Cheng

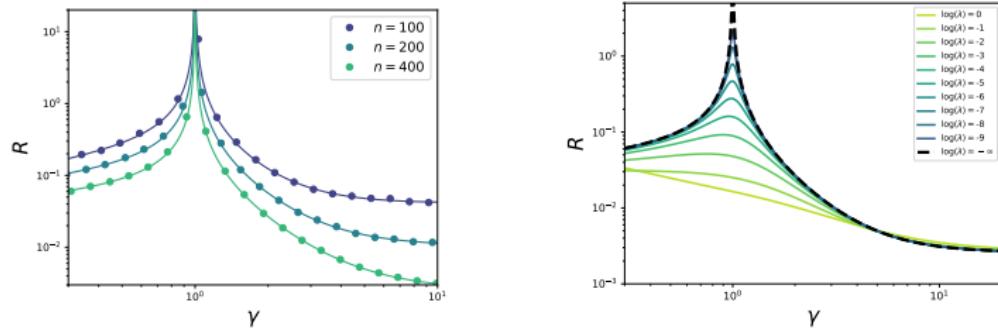
# Surprise #1: N



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

MNIST: 4,000 images 10 classes; 2-layers. Square loss. Belkin, Hsu, Ma, Mandal, 2018

## Surprise #2: Completely general Test error of ridge(less) regression vs $\gamma = p/n$



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

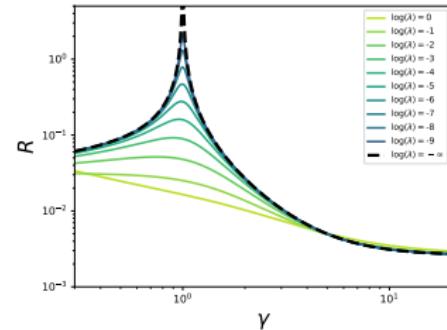
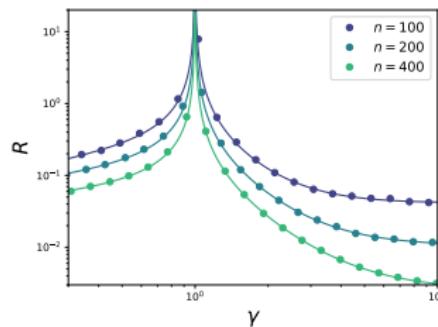
►  $x$ : latent.

$z$ : features

Regress  $y$  vs  $z$

## Surprise #2: Completely general

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



$$y_i = \langle \theta, x_i \rangle + \varepsilon_i, \quad x_i \sim N(\mathbf{0}, I_d),$$

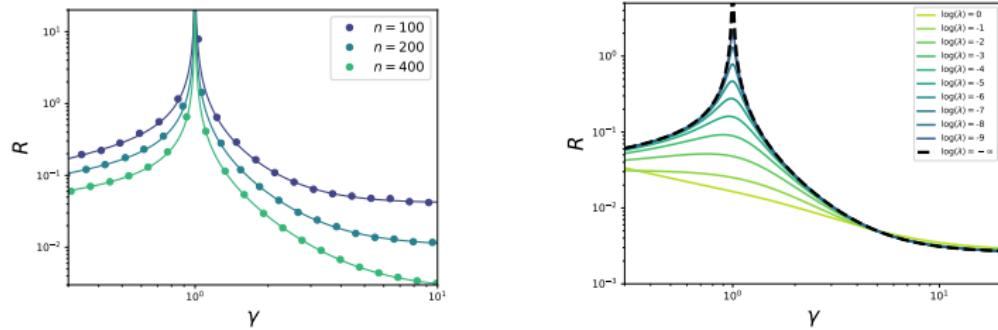
$$z_i = W^T x_i + g_i, \quad W \in \mathbb{R}^{d \times p}, \quad g_i \sim N(\mathbf{0}, I_d).$$

►  $x$ : latent.

$z$ : features

Regress  $y$  vs  $z$

## Surprise #2: Completely general Test error of ridge(less) regression vs $\gamma = p/n$



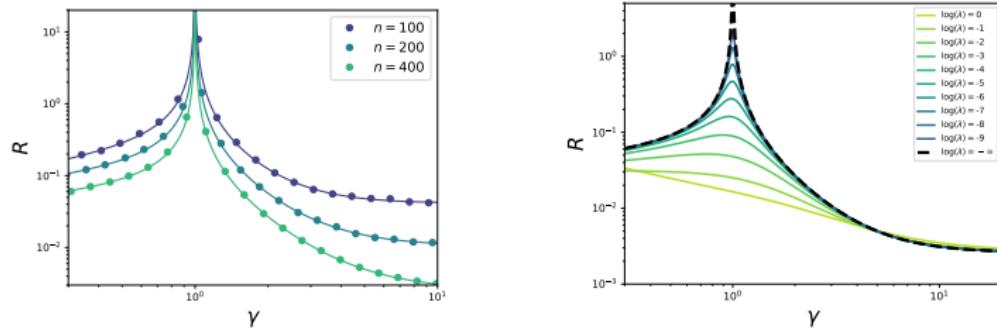
$$y_i = \langle \theta, x_i \rangle + \varepsilon_i, \quad x_i \sim N(0, I_d), \\ z_i = W^T x_i + g_i, \quad W \in \mathbb{R}^{d \times p}, \quad g_i \sim N(0, I_d).$$

►  $x$ : latent.

$z$ : features

Regress  $y$  vs  $z$

## Surprise #2: Completely general



$$y_i = \langle \theta, x_i \rangle + \varepsilon_i, \quad x_i \sim N(0, I_d),$$

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

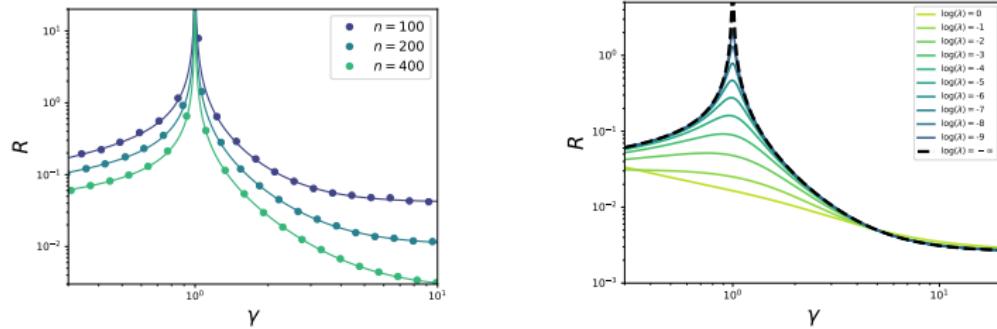
## Equivalent description

$$y_i = \langle \beta, z_i \rangle + \tilde{\varepsilon}_i, \quad z_i \sim N(0, \Sigma_d),$$

$$\Sigma = W W^T + I_p, \quad \beta \in \text{span}(W).$$

General picture?      Connection to neural nets?

## Surprise #2: Completely general



$$y_i = \langle \theta, x_i \rangle + \varepsilon_i, \quad x_i \sim N(0, I_d),$$

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

## Equivalent description

$$y_i = \langle \beta, z_i \rangle + \tilde{\varepsilon}_i, \quad z_i \sim N(0, \Sigma_d),$$

$$\Sigma = W W^T + I_p, \quad \beta \in \text{span}(W).$$

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 Limitations of the linear regime
- 8 Conclusion

## Examples of linear regression

# Setting

## ► Data

$$(y_1, z_1), (y_2, z_2), \dots, (y_n, z_n)$$

$y_i \in \mathbb{R}$ : ‘label’

$z_i \in \mathbb{R}^p$ : ‘features’ vector’.

## ► Distribution

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

---

Potentially  $p = \infty$

# Setting

## ► Data

$$(y_1, z_1), (y_2, z_2), \dots, (y_n, z_n)$$

$y_i \in \mathbb{R}$ : ‘label’

$z_i \in \mathbb{R}^p$ : ‘features’ vector’.

## ► Distribution

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

---

Potentially  $p = \infty$

## Ridge regression

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

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

## Ridge regression

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

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

## Test error

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

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

$$= \|\hat{\beta}(\lambda) - \beta\|_{\Sigma}^2 \quad \Sigma := \mathbb{E}[zz^\top]$$

## Test error

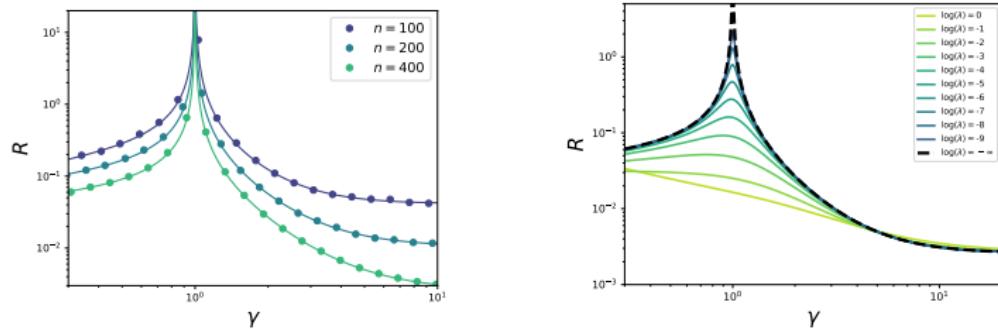
$$\hat{\beta}(\lambda) := \arg \min_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}.$$

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

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

## Examples

# Example #1: Well-concentrated covariates



- ▶  $z_i = \Sigma^{1/2} x_i$ ,  $\mathbb{E}\{x_i x_i^\top\} = I_p$ .
- ▶ Concentration properties for  $x_i$ :

Either: Independent sub-Gaussian coordinates

or: Log-Sobolev

or: ...

## Example #2: Kernel Ridge Regression

### Data

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

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

### Function space view

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

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

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

### Data

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

$$y_i = f_*(x_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{P}), \quad (\text{Hilbert space})$$

### Featurization map

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

Feature space view ( $p = \infty$ )

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

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

---

$$K(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle_{\mathcal{H}_0}.$$

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

### Data

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

$$y_i = f_*(x_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{P}), \quad (\text{Hilbert space})$$

### Featurization map

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

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

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

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

---

$$K(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle_{\mathcal{H}_0}.$$

## 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

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

$$y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in \mathcal{H} \subseteq L^2(\mathbb{P}),$$

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

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

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

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

## Example #4: (Neural) Tangent Regression

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

$$\alpha f(x; \theta^0 + \alpha^{-1} b) = \alpha f(x; \theta^0) + \langle b, \nabla_{\theta} f(x; \theta^0) \rangle + O(\alpha^{-1})$$

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

---

Jacot, Gabriel, Hongler, 2018; Du, Zhai, Poczos, 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

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

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

## Example #4: (Neural) Tangent Regression

### Two-layer neural net

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

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

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

$$\mathbf{z}_i = \Phi_W(\mathbf{x}_i) := (\mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_1^0, \mathbf{x}_i \rangle), \mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_2^0, \mathbf{x}_i \rangle), \dots, \mathbf{x}_i^\top \sigma'(\langle \mathbf{w}_N^0, \mathbf{x}_i \rangle))^\top,$$

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

## Example #4: (Neural) Tangent Regression

Two-layer neural net

$$f(x; \alpha, W) = \sum_{j=1}^N \alpha_j \sigma(\langle w_j, x \rangle)$$

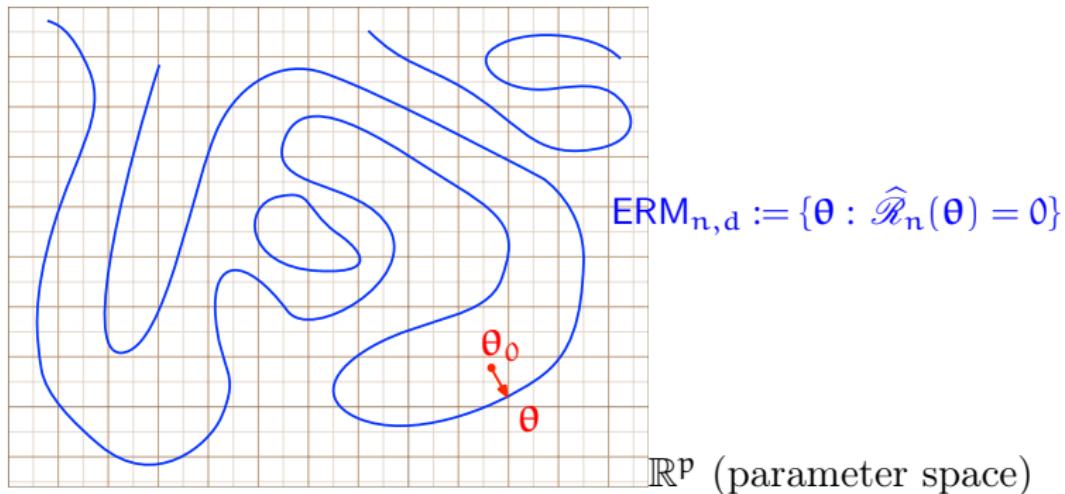
$$f_{NT}(x; \bar{b}, b) = \sum_{j=1}^N \langle b_j, x \rangle \sigma'(\langle w_j^0, x \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle w_j^0, x \rangle).$$

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

$$z_i = \Phi_W(x_i) := (x_i^T \sigma'(\langle w_1^0, x_i \rangle), x_i^T \sigma'(\langle w_2^0, x_i \rangle), \dots, x_i^T \sigma'(\langle w_N^0, x_i \rangle))^T,$$

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

# Heuristic connection with neural nets



## Connection with 'lazy training' of neural nets

'Zero' initialization

$$f(x; w) = \frac{\gamma}{\sqrt{N}} \sum_{j=1}^N b_j \sigma(\langle w_j, x \rangle),$$

$b_1 = \dots = b_{N/2} = +1$ ,  $b_{N/2+1} = \dots = b_N = -1$ , not evolving,

$$(w_j : j \leq N/2) \sim \text{Unif}(\mathbb{S}^{d-1}), \quad w_{N/2+j} = w_j.$$

$$f(x; w) = \frac{\gamma}{\sqrt{N}} \sum_{j=1}^N b_j \sigma(\langle w_j, x \rangle)$$

Theorem (Chizat, Bach, 2018; Oymak, Soltanolkotabi, 2020)

Assume  $\sigma$  is generic,  $x_i \sim N(0, I_d)$ ,  $y_i$  subgaussian. Then for any  $C_0 > 0$  there exists  $C > 0$  such that the following happens. If

$$Nd \geq C d \log d, \quad \gamma \geq C \sqrt{\frac{n^2}{Nd}},$$

then, with probability at least  $1 - 2 \exp(-n/C)$ :

1. Exponential convergence of gradient flow. For all  $t \geq 0$ ,

$$\hat{\mathcal{R}}_n(\theta_t) \leq \hat{\mathcal{R}}_n(\theta_0) e^{-\lambda_* t}, \quad \lambda_* := \gamma^2(d/n)/C$$

2. Neural tangent model is a good approximation

$$\|f(\cdot; \theta_t) - f_{NT}(\cdot; \theta_t^{NT})\|_{L^2(\mathbb{P})} \leq C \left\{ \frac{1}{\gamma} \sqrt{\frac{n^2}{Nd}} + \frac{1}{\gamma^2} \sqrt{\frac{n^5}{Nd^4}} \right\}.$$

$$f(x; w) = \frac{\gamma}{\sqrt{N}} \sum_{j=1}^N b_j \sigma(\langle w_j, x \rangle)$$

Theorem (Chizat, Bach, 2018; Oymak, Soltanolkotabi, 2020)

Assume  $\sigma$  is generic,  $x_i \sim N(0, I_d)$ ,  $y_i$  subgaussian. Then for any  $C_0 > 0$  there exists  $C > 0$  such that the following happens. If

$$Nd \geq C d \log d, \quad \gamma \geq C \sqrt{\frac{n^2}{Nd}},$$

then, with probability at least  $1 - 2 \exp(-n/C)$ :

1. Exponential convergence of gradient flow. For all  $t \geq 0$ ,

$$\hat{\mathcal{R}}_n(\theta_t) \leq \hat{\mathcal{R}}_n(\theta_0) e^{-\lambda_* t}, \quad \lambda_* := \gamma^2(d/n)/C$$

2. Neural tangent model is a good approximation

$$\|f(\cdot; \theta_t) - f_{NT}(\cdot; \theta_t^{NT})\|_{L^2(\mathbb{P})} \leq C \left\{ \frac{1}{\gamma} \sqrt{\frac{n^2}{Nd}} + \frac{1}{\gamma^2} \sqrt{\frac{n^5}{Nd^4}} \right\}.$$

$$f(x; w) = \frac{\gamma}{\sqrt{N}} \sum_{j=1}^N b_j \sigma(\langle w_j, x \rangle)$$

Theorem (Chizat, Bach, 2018; Oymak, Soltanolkotabi, 2020)

Assume  $\sigma$  is generic,  $x_i \sim N(0, I_d)$ ,  $y_i$  subgaussian. Then for any  $C_0 > 0$  there exists  $C > 0$  such that the following happens. If

$$Nd \geq C d \log d, \quad \gamma \geq C \sqrt{\frac{n^2}{Nd}},$$

then, with probability at least  $1 - 2 \exp(-n/C)$ :

1. Exponential convergence of gradient flow. For all  $t \geq 0$ ,

$$\hat{\mathcal{R}}_n(\theta_t) \leq \hat{\mathcal{R}}_n(\theta_0) e^{-\lambda_* t}, \quad \lambda_* := \gamma^2(d/n)/C$$

2. Neural tangent model is a good approximation

$$\|f(\cdot; \theta_t) - f_{NT}(\cdot; \theta_t^{NT})\|_{L^2(\mathbb{P})} \leq C \left\{ \frac{1}{\gamma} \sqrt{\frac{n^2}{Nd}} + \frac{1}{\gamma^2} \sqrt{\frac{n^5}{Nd^4}} \right\}.$$

A general formula

Assumptions:  $p = \infty$  ( $\beta, z_i \in \text{Hilbert}$ )

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

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

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

4.  $u_i := \Sigma^{-1/2} z_i$  satisfies a Hanson-Wright inequality.

Implied by any of the following:

- (a) Independent sub-Gaussian coordinates.
- (b) Concentration 1-Lipschitz convex function

## General result

Theorem (Cheng, M, 2022)

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

Then

$$\mathcal{R}_Z(\lambda) = (1 + \text{err}_n) \mathcal{R}_n^s(\lambda),$$

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

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’  
[‘Historically,’ we proved them independently/earlier]
- ▶  $\mathcal{R}_n^s(\lambda)$  = Risk in a sequence model equivalent

## General result

Theorem (Cheng, M, 2022)

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

Then

$$\mathcal{R}_Z(\lambda) = (1 + \text{err}_n) \mathcal{R}_n^s(\lambda),$$

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

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’  
[‘Historically,’ we proved them independently/earlier]
- ▶  $\mathcal{R}_n^s(\lambda)$  = Risk in a sequence model equivalent

## General result

Theorem (Cheng, M, 2022)

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

Then

$$\mathcal{R}_Z(\lambda) = (1 + \text{err}_n) \mathcal{R}_n^s(\lambda),$$

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

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’  
[‘Historically,’ we proved them independently/earlier]
- ▶  $\mathcal{R}_n^s(\lambda)$  = Risk in a sequence model equivalent

## General result

Theorem (Cheng, M, 2022)

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

Then

$$\mathcal{R}_Z(\lambda) = (1 + \text{err}_n) \mathcal{R}_n^s(\lambda),$$

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

- ▶ Multiplicative error!
- ▶ All that follows will be ‘special cases’  
[‘Historically,’ we proved them independently/earlier]
- ▶  $\mathcal{R}_n^s(\lambda)$  = Risk in a sequence model equivalent

# Key quantity in the sequence model equivalent

**Effective regularization**  $\lambda_*(\lambda)$ :

$$n - \frac{\lambda}{\lambda_*} = \sum_{i \geq 1} \frac{\sigma_i}{\sigma_i + \lambda_*}$$

‘Self-induced’ regularization

$$\lim_{\lambda \rightarrow 0+} \lambda_*(\lambda) = \lambda_*(0) > 0.$$

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

$$\chi_n := \frac{\sigma_{\lfloor \eta n \rfloor} d_\Sigma \log^2(d_\Sigma)}{\kappa n \lambda_*(0)}, \quad d_\Sigma := \max_{k \leq n} \sum_{l \geq k} \frac{\sigma_l}{\sigma_k}.$$

- ▶  $\lambda_* \asymp \sigma_{\lfloor cn \rfloor}$ .
- ▶ Need  $\chi_n \leq n^{1/3-\varepsilon}$
- ▶ Need  $d_\Sigma \leq n^{4/3-\varepsilon}$

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

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

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  $\eta = \eta(C_x) \in (0, 1/2)$  and  $C = C(C_x, D) > 0$  such that the following hold. Define  $\chi_n(\lambda), \kappa, \rho(\lambda)$  as above (with  $\eta = \eta(C_x)$  in Eq. (20)). If it holds that*

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

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

1. **Variance approximation.**

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

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

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

**Remark 9.1** The condition  $|\mathcal{B}| \sim \dots \sim \infty$  in Assumption 1 amounts to requiring that the root

## Equivalent sequence model

$$\theta_i := \langle \beta, v_i \rangle, \quad v_i := i\text{-th eigenvectors of } \Sigma$$

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

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

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

$$\omega^2 = \tau^2 + \mathbb{E}_g \underbrace{\left\{ \sum_{i \geq 1} \sigma_i (\hat{\theta}_i^s - \theta_i)^2 \right\}}_{\mathcal{R}_n^s}, \quad n - \frac{\lambda}{\lambda_*} = \sum_{i \geq 1} \frac{\sigma_i}{\sigma_i + \lambda_*}$$

## Equivalent sequence model

$$\theta_i := \langle \beta, v_i \rangle, \quad v_i \text{ : } i\text{-th eigenvectors of } \Sigma$$

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

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

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

$$\omega^2 = \tau^2 + \mathbb{E}_g \left\{ \underbrace{\sum_{i \geq 1} \sigma_i (\hat{\theta}_i^s - \theta_i)^2}_{\mathcal{R}_n^s} \right\}, \quad n - \frac{\lambda}{\lambda_*} = \sum_{i \geq 1} \frac{\sigma_i}{\sigma_i + \lambda_*}$$

## Specific eigenvalue structures

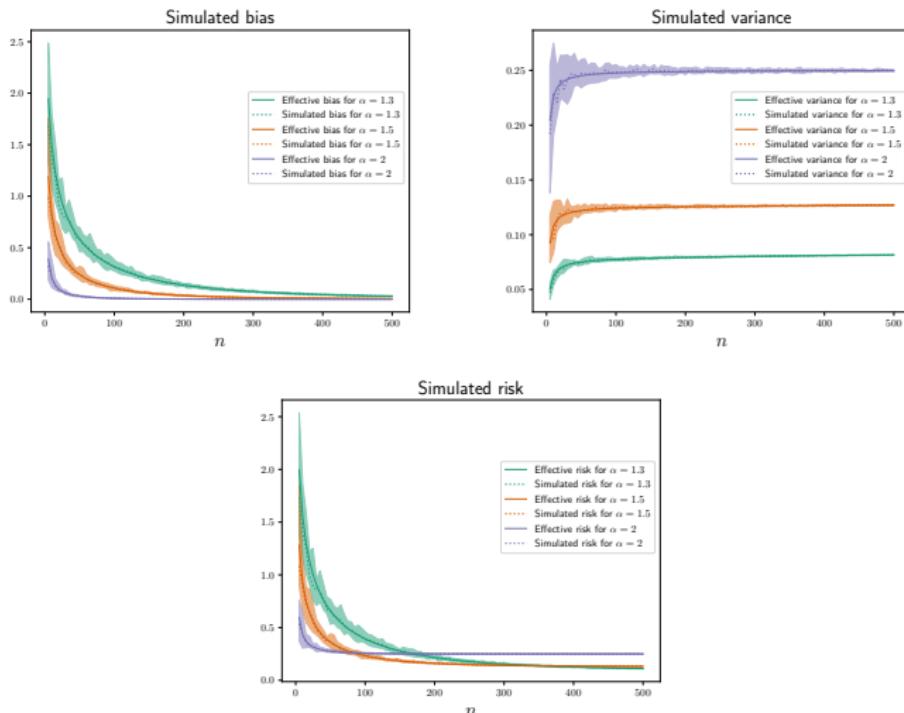
**Power law decay:**  $\sigma_i = i^{-\alpha}, \alpha > 1$

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

---

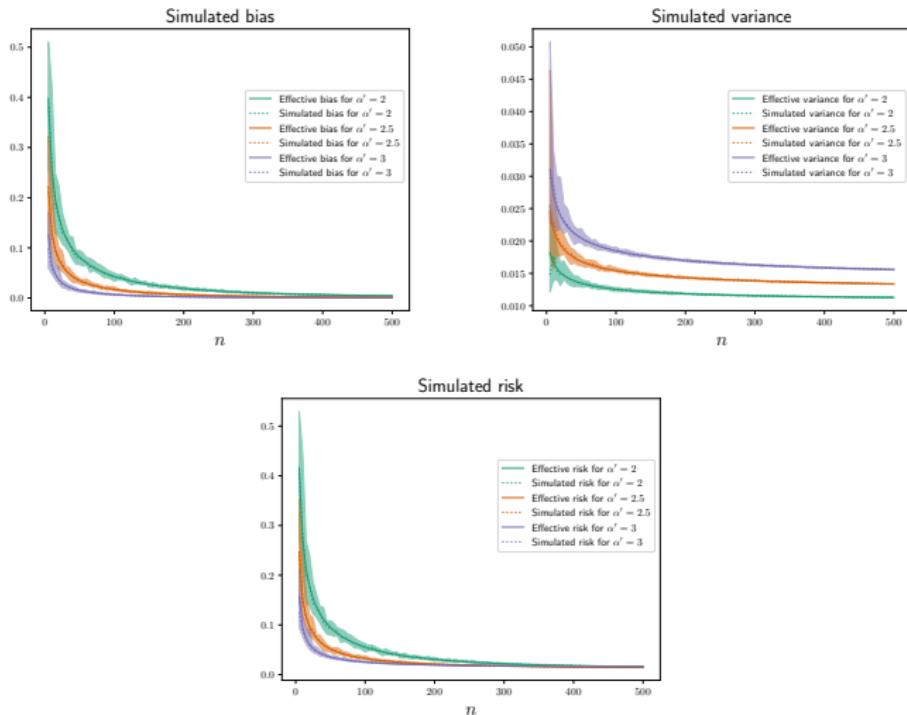
See paper for many other examples

# 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

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

Then  $\mathcal{R}_n^s(\lambda) = B_n^s(\lambda) + V_n^s(\lambda)$ :

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

$$V_n^s(\lambda) = \frac{\tau^2 \text{Tr} \left( \Sigma^2 (\Sigma + \lambda_* \mathbf{I})^{-2} \right)}{n - \text{Tr} \left( \Sigma^2 (\Sigma + \lambda_* \mathbf{I})^{-2} \right)}.$$

## Eigenvalue decay assumption

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

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

Then

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

## Eigenvalue decay assumption

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

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

Then

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

## Eigenvalue decay assumption

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

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

Then

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

## Continuing...

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

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

## Continuing...

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

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

## Continuing...

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

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

Continuing...

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

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

# 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,

$$V_n(\lambda) \leq c_\star \tau^2 \left( \frac{k_\star}{n} + \frac{r_2(k_\star)}{n} \right) \leq c_\star \tau^2 \left( \frac{k_\star}{n} + \frac{4b_{k_\star}^2 n}{\bar{r}(k_\star)} \right),$$

$$B_n(\lambda) \leq c_\star \left( \sigma_{k_\star}^2 \|\beta_{\leq k_\star}\|_{\Sigma^{-1}}^2 + \|\beta_{>k_\star}\|_{\Sigma}^2 \right).$$

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

## Kernel ridge regression

## High-dimensional setting

- ▶  $\mathbf{x}_i \sim \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$
- ▶  $y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad f_* \in L^2(\mathbb{R}^d; \mathbb{P})$
- ▶  $K(x_1, x_2) = h(\langle x_1, x_2 \rangle / d), \quad \mathbb{E}[h(G) H_{\mathcal{K}}(G)] \neq 0$  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:  $\|f_*\|_K \leq C$ , 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:  $\|f_*\|_K \leq C$ , 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

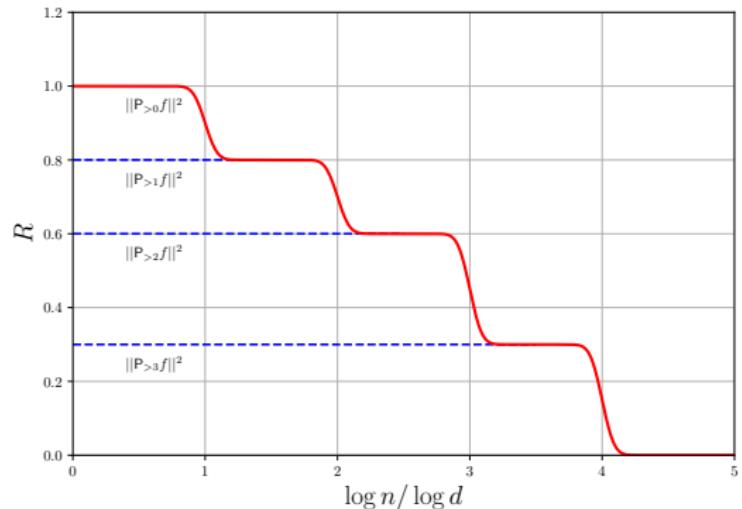
Further, no inner product kernel method can do better.

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

---

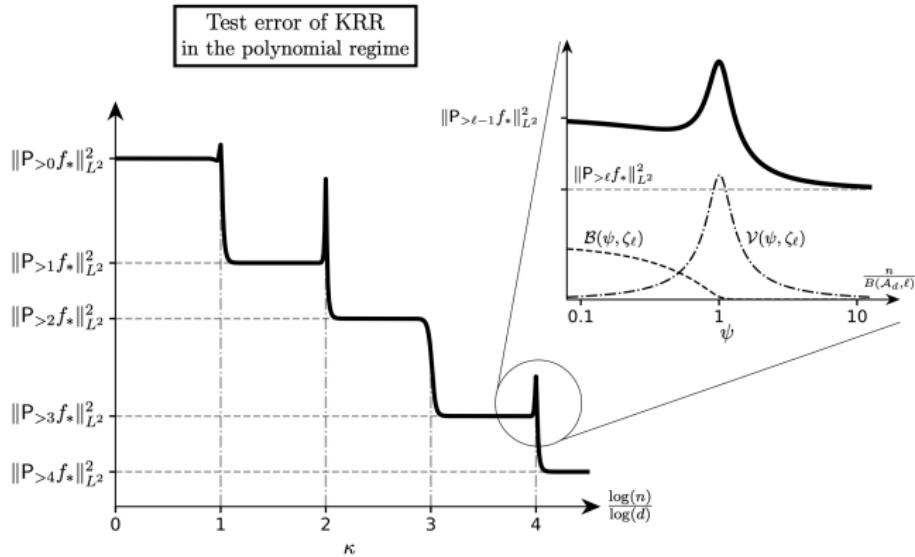
Liang, Rakhlin, Zhai, 2019:  $\|f_*\|_K \leq C$ , upper bounds on the rates. Bartlett, Rakhlin, M., 2021:  
Sharp results for  $n \asymp d$ .

# Sketch



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

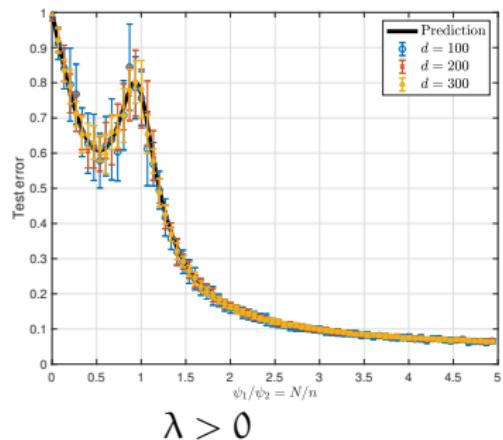
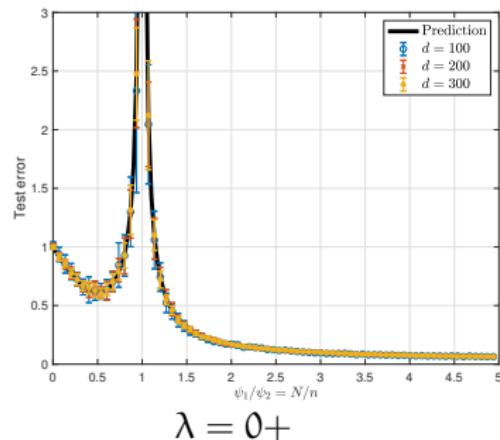
# Developments



Misiakiewicz, 2022; Xiao, Pennington, 2022; Hu, Lu, 2022; **Misiakiewicz, Saeed, 2024**

## Random features

Precise proportional asymptotics:  $N/d \rightarrow \psi_1$ ,  
 $n/d \rightarrow \psi_2$ .



- Solid line: Theoretical prediction

Neural tangent

## Neural Tangent: Parameters view

$$f_{NT}(x; b) = \sum_{j=1}^N \langle b_j, x \rangle \sigma'(\langle w_j, x \rangle).$$

$$z_i := (x_i^T \sigma'(\langle w_1, x_i \rangle), x_i^T \sigma'(\langle w_2, x_i \rangle), \dots, x_i^T \sigma'(\langle w_N, x_i \rangle))^T,$$

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

## Neural Tangent: Function view

$$\hat{f}_{NT} = \operatorname{argmin}_f \left\{ \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{K_N}^2 \right\}.$$

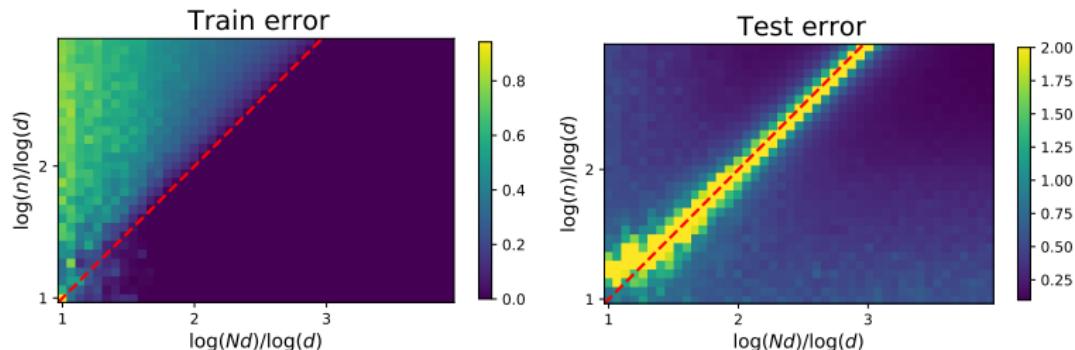
$\|\cdot\|_{K_N}$ : RKHS norm

$$K_N(x_1, x_2) = \frac{1}{Nd} \sum_{i=1}^N \langle x_1, x_2 \rangle \sigma'(\langle w_i, x_1 \rangle) \sigma'(\langle w_i, x_2 \rangle)$$

- ▶ Can we approximate it by  $K(x_1, x_2) = \mathbb{E} K_N(x_1, x_2)$ ?

# A small experiment

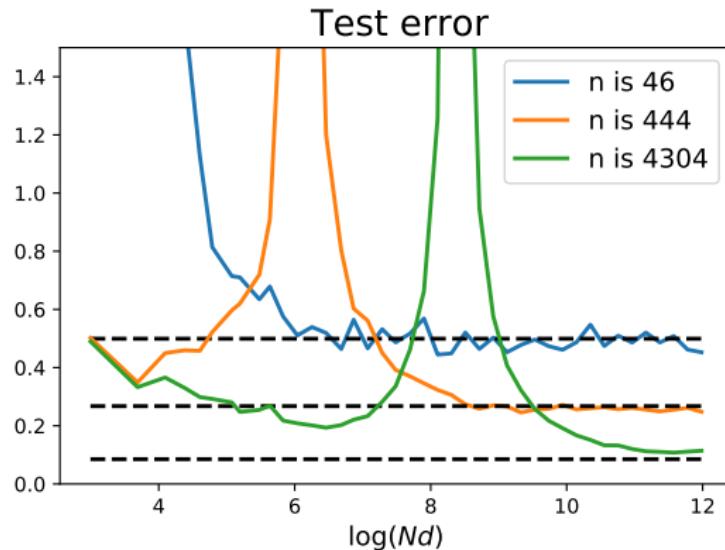
( $d = 20$ )



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

# A small experiment

( $d = 20$ )



- ▶  $f_*(x) = g(\langle \beta_*, x \rangle)$ ,  $g$  = deg-4 polynomial
- ▶ NT ridge regression vs Kernel Ridge(-less) Regression

$$\hat{f} := \arg \min_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \left\{ \|f\|_{\mathcal{K}} \text{ subj. to } f(x_i) = y_i \ \forall i \leq n \right\},$$

## Rigorous confirmation

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

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

Theorem (M, Zhong, 2020, 2021)

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

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

# 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?**

## Limitations of the linear regime

# A simple example

## Data

$$y_i = \sigma(\langle w_*, x_i \rangle) + \varepsilon_i, \quad w_* \in \mathbb{S}^{d-1}, \quad \varepsilon_i \sim N(0, \tau^2).$$

## Machine learning model

$$f(x; W) = \frac{1}{\sqrt{N}} \sum_{j=1}^N a_j \sigma(\langle w_j, x \rangle),$$

$$(a_j^0)_{j \leq N} \sim_{iid} P_A, \quad (w_j^0)_{j \leq N} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$$

What do we know?  $N = 1$

$$f(x; w) = a\sigma(\langle w, x \rangle),$$

$$\widehat{\mathcal{R}}_n(\theta) = \frac{1}{2n} \sum_{i=1}^N (y_i - \gamma a\sigma(\langle w, x \rangle))^2, \quad \theta = (a, w),$$

$$\dot{\theta}_t = -\nabla \widehat{\mathcal{R}}_n(\theta).$$

# What do we know? $N = 1$

$$f(x; w) = \alpha \sigma(\langle w, x \rangle).$$

Proposition (Bai, Mei, M. 2018)

Assume  $\gamma = \alpha = 1$  fixed,  $\sigma \in C^3(\mathbb{R})$ , strictly increasing with bounded derivatives. If  $n \geq Cd \log d$ , then

1. Gradient flow converges exponentially fast:

$$\|w_t - \hat{w}\|_2 \leq C_0 e^{-c_* t}.$$

2. Small generalization error:

$$\mathcal{R}(\hat{w}) - \min_{w \in \mathbb{R}^d} \mathcal{R}(w) \leq C_1 \frac{d \log n}{n}.$$

# What do we know? $N \gg 1$

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

## Proposition

If  $\sigma$  generic,  $N \geq C(\varepsilon) (d \log d \vee n^2/d)$  then:

1. Gradient flow converges to zero risk

$$\widehat{\mathcal{R}}_n(f(\cdot; \theta_t)) \leq \widehat{\mathcal{R}}_n(f(\cdot; \theta_0)) e^{-c_* t}$$

2. The neural tangent model is accurate

$$|\mathcal{R}(f(\cdot; \theta_t)) - \mathcal{R}(f_{NT}(\cdot; \theta_t))| \leq \varepsilon.$$

3. The generalization error is large. For  $d^\ell \ll n \leq d^{\ell+1}$

$$\mathcal{R}(f(\cdot; \hat{\theta})) - \min_{\theta} \mathcal{R}(f(\cdot; \theta)) = \left\| \mathbb{P}_{>\ell} \sigma \right\|_{L^2(N(0,1))}^2 + O(\varepsilon).$$

# What do we know? $N \gg 1$

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

## Proposition

If  $\sigma$  generic,  $N \geq C(\varepsilon) (d \log d \vee n^2/d)$  then:

1. Gradient flow converges to zero risk

$$\widehat{\mathcal{R}}_n(f(\cdot; \theta_t)) \leq \widehat{\mathcal{R}}_n(f(\cdot; \theta_0)) e^{-c_* t}$$

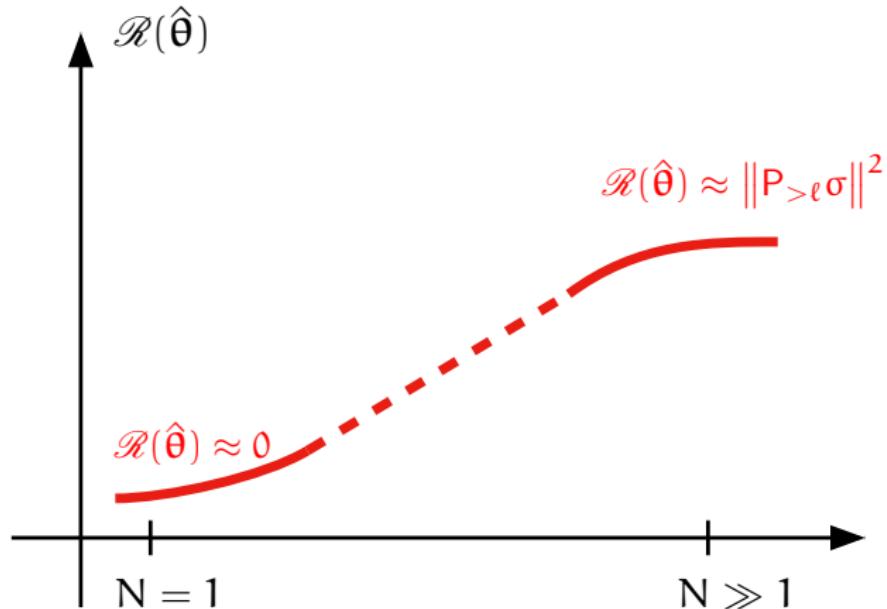
2. The neural tangent model is accurate

$$|\mathcal{R}(f(\cdot; \theta_t)) - \mathcal{R}(f_{NT}(\cdot; \theta_t))| \leq \varepsilon.$$

3. The generalization error is large. For  $d^\ell \ll n \leq d^{\ell+1}$

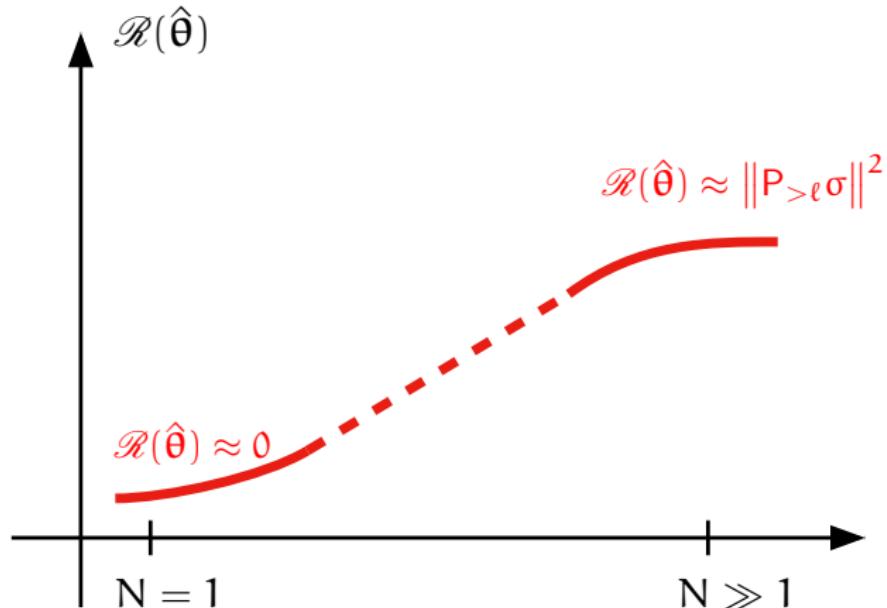
$$\mathcal{R}(f(\cdot; \hat{\theta})) - \min_{\theta} \mathcal{R}(f(\cdot; \theta)) = \left\| P_{>\ell} \sigma \right\|_{L^2(N(0,1))}^2 + O(\varepsilon).$$

## A cartoon



A paradox?

## A cartoon



A paradox?

## Conclusion

# Conclusion

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

Thanks!

# Conclusion

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

Thanks!

# Conclusion

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

Thanks!

# Conclusion

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

Thanks!