



國立臺灣大學  
National Taiwan University

# Neural Tangent Kernel: On Double Descent

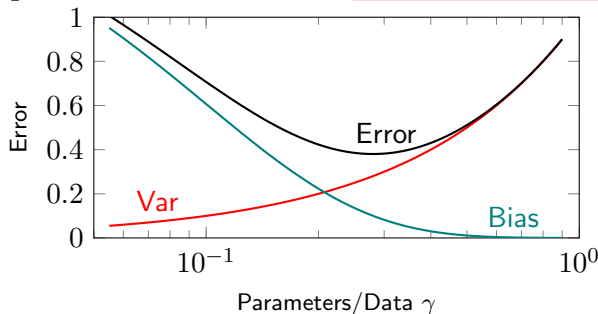
Wen Perng (B10901042)

Physical Theories of (Machine) Learning

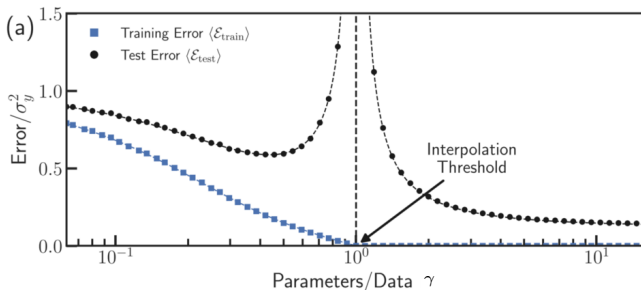
2025 Nov. 12

The generalization error can be separated into bias and variance:

$$\begin{aligned}\mathcal{E}_g &= \mathbb{E}_{\mathbf{x}, \mathcal{D}} \left[ (f_{\star}(\mathbf{x}) - f_{\theta}(\mathbf{x}))^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ \left( \text{Bias}(f_{\theta}, \mathbf{x}) \right)^2 + \text{Var}(f_{\theta}(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ \left( f_{\star}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\theta}(\mathbf{x})] \right)^2 + \mathbb{E}_{\mathcal{D}} \left[ (f_{\theta}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f_{\theta}(\mathbf{x})])^2 \right] \right]\end{aligned}$$



Model performances are not trivial:



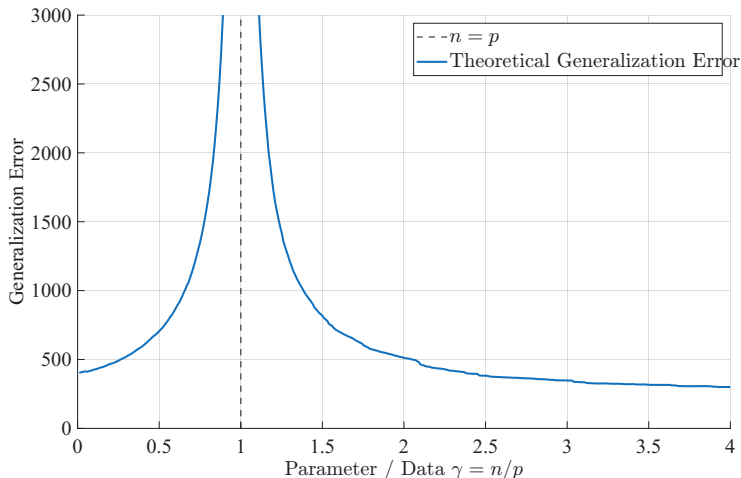
[Rocks et al. '22]

This is the **double descent** behavior.

- Behavior depends on the **hyperparameters**.
- Explain via neural tangent kernel (NTK),  $\Theta_0$ .



Using the **replica trick**, we will be able to derive the following double descent curve for ridgeless kernel regression.





- 1 Review: NTK
- 2 Multiple Descents
- 3 Generalization Error
- 4 Conclusions



- 1 Review: NTK
- 2 Multiple Descents
- 3 Generalization Error
- 4 Conclusions

**Model:** with  $[W_1]_{ij} \sim \mathcal{N}(0, 1)$  and  $[W_2]_{ij} \sim \mathcal{N}(0, \sigma_{W_2}^2)$ ,

$$f_{\theta}(\mathbf{x}) = \frac{W_2}{\sqrt{n_1}} \cdot \sigma \left( \frac{W_1}{\sqrt{n_0}} \mathbf{x} \right). \quad (1)$$

$$1 \underset{1}{\boxed{\phantom{000}}} = 1 \underset{n_1}{\boxed{\phantom{000}}} \cdot \sigma \left( n_1 \underset{n_0}{\boxed{\phantom{000}}} \underset{1}{\boxed{\phantom{00}}} n_0 \right)$$

**Trainable parameters:**

$$\theta = [\text{vec}(W_1); \text{vec}(W_2)] \in \mathbb{R}^{(n_0 n_1 + n_1) \times 1}.$$

**Training data:**  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^p$ , with

$$X = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_p \end{bmatrix} \in \mathbb{R}^{n_0 \times p}, Y = \begin{bmatrix} y_1 & \dots & y_p \end{bmatrix} \in \mathbb{R}^{1 \times p}.$$

NTK describes the learning dynamics under gradient flow:

$$\frac{df_{\theta_t}(\mathbf{x})}{dt} = - (f_{\theta_t}(X) - Y) \Theta_t(X, \mathbf{x}), \quad (2)$$

$$1 \begin{matrix} \boxed{\phantom{00}} \\ 1 \end{matrix} = 1 \begin{matrix} \boxed{\phantom{0000}} \\ p \end{matrix} \begin{matrix} \boxed{\phantom{00}} \\ 1 \end{matrix} \begin{matrix} p \\ \phantom{00} \end{matrix}$$

where

$$\begin{aligned} \Theta_t(\mathbf{x}_i, \mathbf{x}_j) &= \nabla_{\theta}^T f_{\theta_t}(\mathbf{x}_i) \overbrace{\nabla_{\theta} f_{\theta_t}(\mathbf{x}_j)}^{(n_0 n_1 + n_1) \times 1} \\ &= \nabla_{W_1}^T f_{\theta_t}(\mathbf{x}_i) \underbrace{\nabla_{W_1} f_{\theta_t}(\mathbf{x}_j)}_{n_0 n_1 \times 1} + \nabla_{W_2}^T f_{\theta_t}(\mathbf{x}_i) \underbrace{\nabla_{W_2} f_{\theta_t}(\mathbf{x}_j)}_{n_1 \times 1}. \end{aligned} \quad (3)$$





As network width goes to infinity:  $n_0 \rightarrow \infty$  and  $n_1/n_0 = \gamma$  fixed, the neural tangent kernel  $\Theta_t$  becomes **frozen**:

$$\Theta_t \rightarrow \Theta_0 \quad \forall t \geq 0.$$

The final form is a ridgeless kernel regression: denote  $f_{\theta_t}$  by  $f_t$ ,

$$f_{\infty}(\mathbf{x}) = f_0(\mathbf{x}) + \underbrace{(Y - f_0(X)) \Theta_0^{-1}(X, X)}_{\text{regression coefficients}} \underbrace{\Theta_0(X, \mathbf{x})}_{\text{random features}}. \quad (4)$$

As no features are learned, the regression uses the random kernel as random features. This is termed **lazy learning**.



In the wide network limit, the training result is equivalent to a ridgeless ( $\lambda \rightarrow 0$ ) kernel regression with NTK  $\Theta_0$ , with solution of the form<sup>1</sup>

$$f_\infty(\mathbf{x}) = \mathbf{v}^\top \Theta_0(X, \mathbf{x}).$$

By **Mercer's theorem**, the following decomposition exists:

$$\Theta_0(\mathbf{x}, \mathbf{x}') = \sum_{\rho} \psi_{\rho}(\mathbf{x}) \psi_{\rho}(\mathbf{x}'), \quad (5)$$

$$\int p(\mathbf{x}) \Theta_0(\mathbf{x}, \mathbf{x}') \psi_{\rho}(\mathbf{x}) d\mathbf{x} = \eta_{\rho} \psi_{\rho}(\mathbf{x}'), \quad (6)$$

$$\int p(\mathbf{x}) \psi_{\rho}(\mathbf{x}) \psi_{\rho'}(\mathbf{x}) d\mathbf{x} = \eta_{\rho} \delta_{\rho\rho'}. \quad (7)$$

The integrations are over  $\mathbf{x} \in \mathbb{R}^n$ .

---

<sup>1</sup>For simplicity, assume centering  $f_0(X) = 0$ .



By **Moore–Aronszajn theorem**, the NTK forms an RKHS:

$$\mathcal{H} = \left\{ \sum_{\rho} w_{\rho} \psi_{\rho}(\cdot) \right\} = \left\{ \mathbf{w}^{\top} \boldsymbol{\psi}(\cdot) \right\} \quad (8)$$

with

$$\Lambda := \mathbb{E}_{\mathbf{x}} \left[ \boldsymbol{\psi}(\mathbf{x}) \boldsymbol{\psi}^{\top}(\mathbf{x}) \right] = \text{diag}(\eta_1, \eta_2, \dots). \quad (9)$$

The teacher (target) and the student (learned) functions both lies in  $\mathcal{H}$ :<sup>2</sup>

$$f_{\star}(\cdot) = \mathbf{w}_{\star}^{\top} \boldsymbol{\psi}(\cdot),$$

$$f_{\boldsymbol{\theta}}(\cdot) = \mathbf{w}^{\top} \boldsymbol{\psi}(\cdot).$$

depends on training

<sup>2</sup>If  $f_{\star} \notin \mathcal{H}$ , it is equivalent to having noise in the sampling of  $f_{\star}$ .



- 1 Review: NTK
- 2 Multiple Descents
- 3 Generalization Error
- 4 Conclusions



Let  $\mathbf{x}$  follow the same data distribution as the training dataset  $\mathcal{D}$ ,

$$\begin{aligned}\mathcal{E}_g &= \mathbb{E}_{\mathbf{x}} \left[ \|f_{\star}(\mathbf{x}) - f_{\infty}(\mathbf{x})\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ (f_{\star}(\mathbf{x}) - f_0(\mathbf{x}))^2 \right] \\ &\quad - 2\mathbb{E}_{\mathbf{x}} \left[ (f_{\star}(\mathbf{x}) - f_0(\mathbf{x})) (Y - f_0(X)) \Theta_0^{-1}(X, X) \Theta(X, \mathbf{x}) \right] \\ &\quad + \mathbb{E}_{\mathbf{x}} \left[ ((Y - f_0(X)) \Theta_0^{-1}(X, X) \Theta_0(X, \mathbf{x}))^2 \right] \\ &= a + \text{Tr} \left\{ B \Theta_0^{-1}(X, X) \right\} \\ &\quad + \sum c_i \text{Tr} \left\{ C_i \Theta_0^{-1}(X, X) D_i \Theta_0^{-1}(X, X) \right\}\end{aligned}$$

The generalization depends on the **spectrum** of the **empirical** NTK.

For the two-layer network:

$$f_{\theta}(\mathbf{x}) = \frac{W_2}{\sqrt{n_1}} \cdot \sigma \left( \frac{W_1}{\sqrt{n_0}} \mathbf{x} \right) \quad (1)$$

$$\nabla_{W_1} f_{\theta}(\mathbf{x}) \in \mathbb{R}^{n_0 n_1 \times 1}, \nabla_{W_2} f_{\theta}(\mathbf{x}) \in \mathbb{R}^{n_1 \times 1}.$$

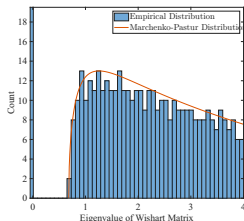
Thus, the empirical NTK is

$$\begin{aligned} \Theta_0(X, X) = & \underbrace{p}_{n_0 n_1} \underbrace{\begin{bmatrix} \vdots \\ \nabla_{W_1}^{\top} f_{\theta}(\mathbf{x}_i) \\ \vdots \end{bmatrix}}_{n_0 n_1} \underbrace{\begin{bmatrix} \cdots \nabla_{W_1} f_{\theta}(\mathbf{x}_j) \cdots \end{bmatrix}}_p \underbrace{n_0 n_1}_{n_0 n_1} \\ & + \underbrace{p}_{n_1} \underbrace{\begin{bmatrix} \vdots \\ \nabla_{W_2}^{\top} f_{\theta}(\mathbf{x}_i) \\ \vdots \end{bmatrix}}_{n_1} \underbrace{\begin{bmatrix} \cdots \nabla_{W_2} f_{\theta}(\mathbf{x}_j) \cdots \end{bmatrix}}_p \underbrace{n_1}_{n_1}. \end{aligned}$$

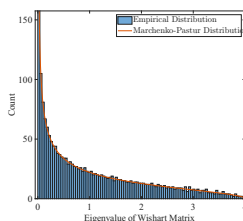
The empirical NTK contains terms of the form

$$\frac{1}{n} F^{\top} F,$$

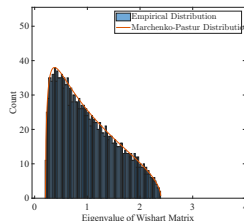
where  $F$  is of size  $n_0 n_1 \times p$  or  $n_1 \times p$  with **stochastic** entries. The limiting spectrum is of **Marchenko–Pastur**-like form<sup>3</sup>.



(a)  $n/p = 0.67$



(b)  $n/p = 1.00$

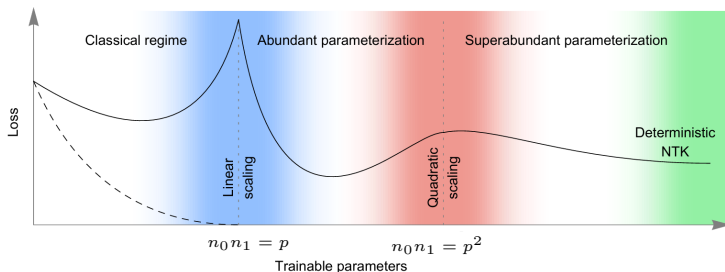


(c)  $n/p = 3.33$

<sup>3</sup>The figure above is only an illustration if there is no nonlinearity.

The rank of NTK undergoes **transition** at  $n_0 n_1 = p$  and  $n_1 = p$ .

The figure below shows an illustration of  $n_0 = n_1 \rightarrow \infty$ .



[Adlam et al. '20]

A quantitative analysis requires knowledge in **nonlinear random matrix theory** [AP20; S23].





- 1 Review: NTK
- 2 Multiple Descents
- 3 Generalization Error**
- 4 Conclusions



Since NTK training is similar to **kernel ridge regression**, let us study the generalization behavior of the latter first!

---

There is label noise:  $y_i = f_\star(\mathbf{x}_i) + \varepsilon_i$  and  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . The dataset is identified as  $\mathcal{D} = \{(\mathbf{x}_i, \varepsilon_i)\}_{i=1}^p$ .

With these, the training loss function for kernel ridge regression is

$$\begin{aligned}\mathcal{E}_\lambda(\mathbf{w}) &= \sum_{i \in \mathcal{D}} (y_i - f_\theta(\mathbf{x}_i))^2 + \lambda \|\mathbf{f}_\theta\|_{\mathcal{H}}^2 \\ &= \sum_{i \in \mathcal{D}} \left( (\mathbf{w}_\star - \mathbf{w})^\top \boldsymbol{\psi}(\mathbf{x}_i) + \varepsilon_i \right)^2 + \lambda \|\mathbf{w}\|^2 \\ &= \sum_{i \in \mathcal{D}} \left( \bar{\mathbf{w}}^\top \boldsymbol{\psi}(\mathbf{x}_i) + \varepsilon_i \right)^2 + \lambda \|\mathbf{w}_\star - \bar{\mathbf{w}}\|^2,\end{aligned}$$

where  $\bar{\mathbf{w}} = \mathbf{w}_\star - \mathbf{w}$ . Denote  $\bar{\mathbf{w}}^* := \arg \min \mathcal{E}_\lambda(\bar{\mathbf{w}})$ .



Similarly, the generalization error will be:

$$\begin{aligned}\mathcal{E}_g(\overline{\mathbf{w}}^*) &= \mathbb{E}_{\mathbf{x}} \left[ (f_{\star}(\mathbf{x}) - f_{\boldsymbol{\theta}}(\mathbf{x}))^2 \right] \\ &= (\mathbf{w}_{\star} - \mathbf{w}^*)^{\top} \mathbb{E}_{\mathbf{x}} \left[ \boldsymbol{\psi}(\mathbf{x}) \boldsymbol{\psi}^{\top}(\mathbf{x}) \right] (\mathbf{w}_{\star} - \mathbf{w}^*) \\ &= \overline{\mathbf{w}}^{*\top} \Lambda \overline{\mathbf{w}}^*.\end{aligned}$$

The typical generalization error will be

$$\begin{aligned}\mathcal{E}_g &= \mathbb{E}_{\mathcal{D}} [\mathcal{E}_g(\overline{\mathbf{w}}^*)] \\ &= \mathbb{E}_{\mathcal{D}} \left[ \int \mathcal{E}_g(\overline{\mathbf{w}}) \delta(\overline{\mathbf{w}} - \overline{\mathbf{w}}^*) \, \mathrm{d}\overline{\mathbf{w}} \right].\end{aligned}$$



$$\mathcal{E}_g = \mathbb{E}_{\mathcal{D}} \left[ \int \mathcal{E}_g(\overline{\mathbf{w}}) \delta(\overline{\mathbf{w}} - \overline{\mathbf{w}}^*) d\overline{\mathbf{w}} \right]$$

Consider

$$Z = \int d\overline{\mathbf{w}} \exp \left\{ \beta \left( -\mathcal{E}_\lambda(\overline{\mathbf{w}}) \right) \right\},$$

then

$$\delta(\overline{\mathbf{w}} - \overline{\mathbf{w}}^*) = \lim_{\beta \rightarrow \infty} \frac{1}{Z} e^{-\beta \mathcal{E}_\lambda(\overline{\mathbf{w}})}.$$

How do we incorporate the  $\mathcal{E}_g$  term? By including a **source term**.

By utilizing the partition function [CBP21]:

$$Z = \int d\bar{\mathbf{w}} \exp \left\{ \beta \left( -\mathcal{E}_\lambda(\bar{\mathbf{w}}) + t \cdot \mathcal{E}_g(\bar{\mathbf{w}}) \right) \right\}, \quad (10)$$

the test error under kernel ridge regression is

$$\mathcal{E}_g(\bar{\mathbf{w}}^*) = \lim_{\beta \rightarrow \infty} \frac{1}{\beta} \frac{\partial}{\partial t} \log Z \Big|_{t=0}. \quad (11)$$

$$\bar{\mathbf{w}}^* = \arg \min \mathcal{E}_\lambda(\bar{\mathbf{w}})$$

Then, the generalization error is

$$\mathcal{E}_g = \mathbb{E}_{\mathcal{D}} [\mathcal{E}_g(\bar{\mathbf{w}}^*)] = \lim_{\beta \rightarrow \infty} \frac{1}{\beta} \frac{\partial}{\partial t} \mathbb{E}_{\mathcal{D}} [\log Z] \Big|_{t=0}. \quad (12)$$

Thus, it requires the **replica trick**<sup>4</sup>!

<sup>4</sup>I apologize for using a different notation as in the paper.



Recap: what is the **replica trick**?

$$\mathbb{E}_{\mathcal{D}}[\log Z] = \mathbb{E}_{\mathcal{D}} \left[ \lim_{x \rightarrow \infty} \frac{Z^x - 1}{x} \right] = \lim_{x \rightarrow 0} \frac{\mathbb{E}_{\mathcal{D}}[Z^x] - 1}{x}.$$

Steps:

1. For integer  $n$ , compute  $Z^n$  with replicated state parameter  $w^a$  and identical disorder parameter  $\mathcal{D}$ .
2. Apply  $\mathbb{E}_{\mathcal{D}}[\cdot]$ .
3. Along the way, apply
  - ▶ Hubbard–Stratonovich transformation, and
  - ▶ Fourier transform relation of delta distribution.
4. Compute saddle point equation.

23/37

$$\mathbb{E}_{\mathcal{D}}[Z^n] = \int \mathcal{D}[\bar{\mathbf{w}}] \cdot e^{-\beta \sum_{a=1}^n (\lambda \|\mathbf{w}_* - \bar{\mathbf{w}}^a\|^2 - t \mathcal{E}_g(\bar{\mathbf{w}}^a))} e^{-\frac{p}{2} \log \det(1 + 2\beta C)},$$

where

$$\begin{aligned} \blacksquare &= \int \mathcal{D}[C] e^{-\frac{p}{2} \log \det(1 + 2\beta C)} \prod_{a \geq a'} \delta(C^{ab} - \bar{\mathbf{w}}^a{}^\top \Lambda \bar{\mathbf{w}}^b - \sigma^2) \\ &\quad \underbrace{\int \frac{d\hat{C}^{ab}}{\sqrt{2\pi}} e^{i\hat{C}^{ab}(C^{ab} - \bar{\mathbf{w}}^a{}^\top \Lambda \bar{\mathbf{w}}^b - \sigma^2)}}_{\text{Fourier representation of the delta function}} \\ &= \int \frac{\mathcal{D}[C] \mathcal{D}[\hat{C}]}{(2\pi)^{\frac{n(n+1)}{2}}} e^{-\frac{p}{2} \log \det(1 + 2\beta C) + i \sum_{a \geq b} \hat{C}^{ab} (C^{ab} - \bar{\mathbf{w}}^a{}^\top \Lambda \bar{\mathbf{w}}^b - \sigma^2)}, \end{aligned}$$

with  $\mathcal{D}[C] = \prod_{a \geq b} dC^{ab}$  and  $\mathcal{D}[\hat{C}] = \prod_{a \geq b} d\hat{C}^{ab}$ .





$$\begin{aligned}
 \mathbb{E}_{\mathcal{D}}[Z^n] &= \int \frac{\mathcal{D}[C] \mathcal{D}[\hat{C}] \mathcal{D}[\bar{\mathbf{w}}]}{(2\pi)^{\frac{n(n+1)}{2}}} e^{-\beta \sum_a (\lambda \|\mathbf{w}_\star - \bar{\mathbf{w}}^a\|^2 - t \bar{\mathbf{w}}^a \mathbf{T} \Lambda \bar{\mathbf{w}}^a)} \\
 &\quad \cdot e^{-\frac{p}{2} \log \det(\mathbf{1} + 2\beta C) + i \sum_{a \geq b} \hat{C}^{ab} (C^{ab} - \bar{\mathbf{w}}^a \mathbf{T} \Lambda \bar{\mathbf{w}}^b - \sigma^2)} \\
 &= e^{-\frac{n(n+1)}{2} \log(2\pi) - \beta \lambda n \mathbf{w}_\star \mathbf{T} \mathbf{w}_\star} \\
 &\quad \times \int \mathcal{D}[C] \mathcal{D}[\hat{C}] e^{-p G_E} e^{-G_S} e^{i \sum_{a \geq b} \hat{C}^{ab} (C^{ab} - \sigma^2)}
 \end{aligned}$$

Note that  $e^{-G_S}$  is a Gaussian integral in  $\bar{\mathbf{w}}$ :

$$e^{-G_S} = \int d\bar{\mathbf{w}} \cdot \exp \left\{ -\beta \bar{\mathbf{w}} \mathbf{T} X \bar{\mathbf{w}} + \beta \boldsymbol{\xi} \mathbf{T} \bar{\mathbf{w}} \right\}.$$

$$X^{ab} = (\lambda \mathbf{1} - t \Lambda + \frac{i}{2\beta} \hat{C}^{aa} \Lambda) \delta_{ab} + \frac{i}{2\beta} \hat{C}^{ab} \Lambda$$

$$\boldsymbol{\xi}^a = 2\lambda \mathbf{w}_\star$$

$$[\bar{\mathbf{w}}^1; \dots; \bar{\mathbf{w}}^n]$$

We thus have

$$e^{-G_S} = (\det(\beta X/\pi))^{-1/2} e^{\frac{1}{4\beta} \boldsymbol{\xi}^\top X^{-1} \boldsymbol{\xi}}.$$

Combining all the above, we obtain the integral

$$\mathbb{E}_{\mathcal{D}} [Z^n] = e^{O(n)} \int \mathcal{D}[C] \mathcal{D}[\hat{C}] e^{-n\beta S[C, \hat{C}]},$$

$$\begin{aligned} S[C, \hat{C}] = & \frac{p}{2n\beta} \log \det(\mathbb{1} + 2\beta C) + \frac{1}{2n\beta} \log \det X \\ & + \frac{1}{4n\beta^2} \boldsymbol{\xi}^\top X^{-1} \boldsymbol{\xi} - \frac{i}{n\beta} \sum_{a \geq b} \hat{C}^{ab} (C^{ab} - \sigma^2). \end{aligned}$$

This requires **saddle-point method** to find the typical  $C$  and  $\hat{C}$ !

By the **convexity** of the problem, we can assume the symmetric replica ansatz:

$$C_0 = C^{aa}, \quad \hat{C}_0 = \hat{C}^{aa}, \quad C_1 = C^{a \neq b}, \quad \hat{C}_1 = \hat{C}^{a \neq b}.$$

Knowing that  $n \rightarrow 0$  and  $p \rightarrow \infty$ ,

$$\begin{aligned} S[C, \hat{C}] = & \left[ \frac{p}{2\beta} \left( \log(1 + 2\beta(C_0 - C_1)) + \frac{2\beta C_1}{1 + 2\beta(C_0 - C_1)} \right) \right. \\ & + \frac{1}{2\beta} \left( \log \det \left( \lambda \mathbb{1} + \left( \frac{i}{\beta} \left( \hat{C}_0 - \frac{1}{2} \hat{C}_1 \right) - t \right) \Lambda \right) + \text{Tr} \left\{ \frac{i}{2\beta} \hat{C}_1 \Lambda \left( \lambda \mathbb{1} + \left( \frac{i}{\beta} \left( \hat{C}_0 - \frac{1}{2} \hat{C}_1 \right) - t \right) \Lambda \right)^{-1} \right\} \right) \\ & + \frac{1}{4\beta^2} \mathbf{w}_*^T \left( \lambda \mathbb{1} + \left( \frac{i}{\beta} \left( \hat{C}_0 - \frac{1}{2} \hat{C}_1 \right) - t \right) \Lambda \right)^{-1} \mathbf{w}_* \\ & \left. - \frac{i}{\beta} \left( \hat{C}_0(C_0 - \sigma^2) - \frac{1}{2} \hat{C}_1(C_1 - \sigma^2) \right) \right] + O(n). \end{aligned}$$

The saddle point  $(C^*, \hat{C}^*)$  then satisfies

$$0 = \frac{\partial S}{\partial C_0} = \frac{\partial S}{\partial C_1} = \frac{\partial S}{\partial \hat{C}_0} = \frac{\partial S}{\partial \hat{C}_1}.$$

By defining

$$\kappa(t) := \lambda (1 + 2\beta(C_0^* - C_1^*)),$$

one obtains

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[Z^n] &\approx e^{O(n)} e^{-n\beta S[C^*, \hat{C}^*]} \\ \Rightarrow \mathbb{E}_{\mathcal{D}}[\log Z] &= \lim_{n \rightarrow 0} \frac{1}{n} \left( O(n) - n\beta S[C^*, \hat{C}^*] \right) \\ \Rightarrow \mathcal{E}_g &= \lim_{\beta \rightarrow \infty} \lim_{n \rightarrow 0} \left. \frac{\partial}{\partial t} S[C^*, \hat{C}^*] \right|_{t=0} \\ &\stackrel{!}{=} \frac{\zeta}{1 - \zeta} \left( \sigma^2 + \kappa^2 \sum_{\rho} \frac{\eta_{\rho} w_{\star, \rho}^2}{(p\eta_{\rho} + \kappa)^2} \right) + \kappa^2 \sum_{\rho} \frac{\eta_{\rho} w_{\star, \rho}^2}{(p\eta_{\rho} + \kappa)^2},\end{aligned}$$

where  $\kappa := \kappa(0)$  and  $\kappa'(0) = \kappa^2 \zeta / (1 - \zeta)$ .<sup>5</sup>

---

<sup>5</sup>I have checked till the last step.

$$\mathcal{E}_g = \underbrace{\frac{\zeta}{1-\zeta} \left( \sigma^2 + \kappa^2 \sum_{\rho} \frac{\eta_{\rho} w_{\star, \rho}^2}{(p\eta_{\rho} + \kappa)^2} \right)}_{\text{variance}} + \underbrace{\kappa^2 \sum_{\rho} \frac{\eta_{\rho} w_{\star, \rho}^2}{(p\eta_{\rho} + \kappa)^2}}_{\text{bias}}, \quad (13)$$

where<sup>6</sup>

$$\kappa := \kappa(0) = \lambda + \sum_{\rho} \frac{\kappa \eta_{\rho}}{p\eta_{\rho} + \kappa}, \quad \zeta = \sum_{\rho} \frac{p\eta_{\rho}^2}{(\kappa + p\eta_{\rho})^2}. \quad (14)$$

Note: The function  $\kappa(t)$  actually has close relation with the spectrum of the empirical NTK.

<sup>6</sup>Note that if  $f_{\star} \notin \mathcal{H}$ , some additional terms need to be included:  $\sigma^2 \mapsto \sigma^2 + \|\mathbf{a}_{\star}\|^2$  and  $\mathcal{E}_g \mapsto \mathcal{E}_g + \|\mathbf{a}_{\star}\|^2$ , where  $\mathbf{a}_{\star}$  is the out-of-RKHS part of  $f_{\star}$ . See the paper [CBP21] for the full picture.



Consider a linear regression with  $\psi_\rho : \mathbb{R}^D \rightarrow \mathbb{R}$ ,  $\mathbf{x} \mapsto x_\rho$ ,  $\eta_\rho = 1$ ,  $D = 400$ ,  $p = 100$ ,  $\sigma = 0.5$ , and  $\mathbf{w}_\star$ ,  $\mathbf{x}$ 's are all normally distributed. Setting  $\rho = 1, \dots, n$ , we have:

$$\kappa = \cancel{\lambda}^0 + \sum_{\rho=1}^n \frac{\cancel{\kappa\eta_\rho}^1}{\cancel{p\eta_\rho}^1 + \kappa} = \frac{n\kappa}{p + \kappa},$$

$$\kappa^2 - (n - p)\kappa = 0 \rightsquigarrow \kappa = \begin{cases} n - p, & n \geq p; \\ 0, & n < p. \end{cases}$$

Similarly,

$$\zeta = \begin{cases} p/n, & n \geq p; \\ n/p, & n < p. \end{cases}$$



Denote  $\gamma = n/p = \text{parameter} / \text{data}$ ,

$$\|\mathbf{w}_{\star,n}\|^2 = \sum_{\rho=1}^n w_{\star,\rho}^2 \text{ and } \|\mathbf{a}_{\star,n}\|^2 = \sum_{\rho=n+1}^{\infty} w_{\star,\rho}^2.$$

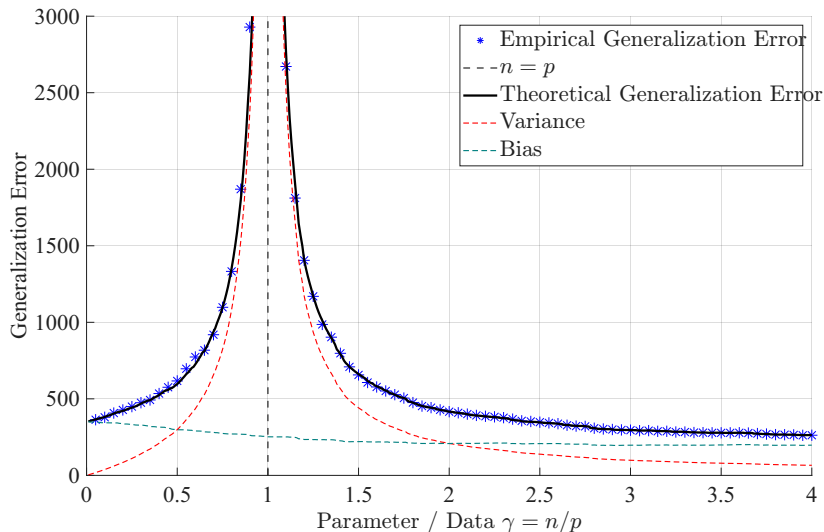
This example has closed form solution for the generalization error:

$$\mathcal{E}_g = \begin{cases} \frac{\gamma-1}{\gamma} \|\mathbf{w}_{\star,n}\|^2 + \frac{\sigma^2 + \|\mathbf{a}_{\star,n}\|^2}{\gamma-1} + \|\mathbf{a}_{\star,n}\|^2, & \gamma \geq 1; \\ \frac{\gamma}{1-\gamma} (\sigma^2 + \|\mathbf{a}_{\star,n}\|^2) + \|\mathbf{a}_{\star,n}\|^2, & 1 > \gamma \geq 0. \end{cases}$$

## Example – Double Descent (3)



Run the code yourself! [[link](#)]

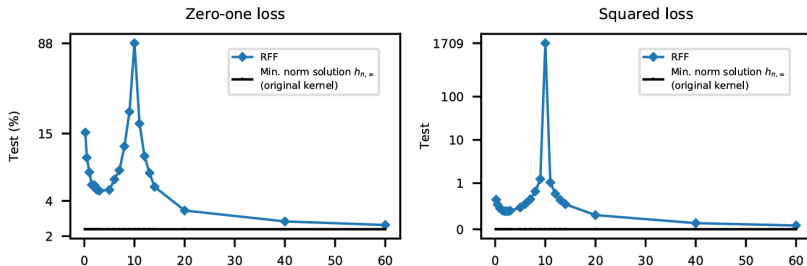






- 1 Review: NTK
- 2 Multiple Descents
- 3 Generalization Error
- 4 Conclusions

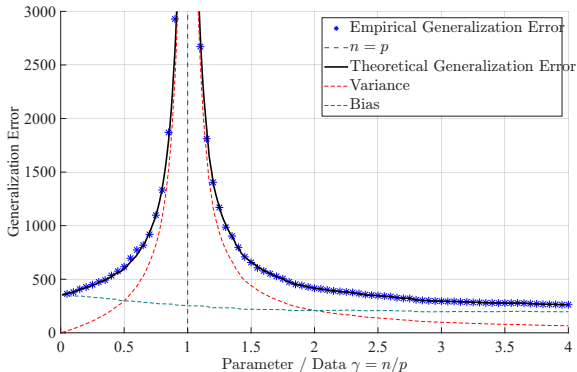
The double descent phenomena have been seen not only in regression or toy model, but also in real dataset classification (e.g., MNIST).



[Belkin et al. '19]

It is known that for neural networks (MLP) with large width, its behavior is similar to ridgeless kernel regression. This talk then linked it with the double descent in kernel regression.

The blowup is caused by variance of small-eigenvalue features from the empirical NTK.



Using the replica trick, we saw the non-monotonic behaviors come from the variance.



- [AP20] Ben Adlam and Jeffrey Pennington. *The Neural Tangent Kernel in High Dimensions: Triple Descent and a Multi-Scale Theory of Generalization*. 2020. arXiv: 2008.06786 [stat.ML]. URL: <https://arxiv.org/abs/2008.06786>.
- [S23] Roland Speicher. *High Dimensional Analysis: Random Matrices and Machine Learning*. Lecture at Saarland University. 2023. URL: <https://rolandspeicher.com/lectures/course-on-high-dimensional-analysis-random-matrices-and-machine-learning-summer-term-2023/>.



- [CBP21] Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. “Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks”. In: *Nature Communications* 12.1 (May 2021). ISSN: 2041-1723. DOI: [10.1038/s41467-021-23103-1](https://doi.org/10.1038/s41467-021-23103-1). URL: <http://dx.doi.org/10.1038/s41467-021-23103-1>.