

# Exact Gap between Generalization Error and Uniform Convergence in Random Feature Models(2021, arxiv)

---

SeongSik Choi

January 26, 2022

Seoul National University

# Motivation

- 학습 이론에서, test err와 train err의 차이가 작을수록 generalize 성능이 좋다는 의미.
- Uniform convergence(UC) is the supremum of (test err - train err) over a certain function class.
- $\therefore$  UC is a powerful tool for understanding the generalization performance of predictors.
- Bounds on uniform convergence는 일반적으로  $\sqrt{\text{complexity}/n}$ 의 형태이므로 만약 tight하다면 function class가 복잡할수록 generalize를 못 해야한다. (그러나 현실은 DNN도 generalize를 잘 한다.)
- **Question** : What is the exact gap between uniform convergence and the actual generalization errors?

# Model setup

- Consider a dataset  $(\mathbf{x}_i, y_i)_{i \in [n]}$  with  $n$  samples.
- The covariates follow  $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ .
- Responses satisfy  $y_i = f_d(\mathbf{x}_i) + \varepsilon_i$ .
- The noises satisfying  $\varepsilon_i \sim_{iid} \mathcal{N}(0, \tau^2)$  indep of  $(\mathbf{x}_i)_{i \in [n]}$ .
- The noisy setting ( $\tau^2 > 0$ ) and noiseless setting ( $\tau^2 = 0$ ).
- Function class : Random features(RF) function class.

## Random features model

**(Random features model)** Let  $(\boldsymbol{\theta}_j)_{j \in [M]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$  be the random feature vectors. Given an activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , we define the random features function class  $\mathcal{F}_{\text{RF}}(\boldsymbol{\Theta})$  by

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\Theta}) \equiv \left\{ f_{\mathbf{a}}(\mathbf{x}) = \sum_{j=1}^M a_j \sigma \left( \langle \mathbf{x}, \boldsymbol{\theta}_j \rangle / \sqrt{d} \right) : \mathbf{a} \in \mathbb{R}^M \right\}$$

$$\min_{f_{\mathbf{a}} \in \mathcal{F}_{\text{RF}}(\boldsymbol{\Theta})} R(f_{\mathbf{a}}) = \min_{\mathbf{a} \in \mathbb{R}^M} R(\mathbf{a})$$

This function class is proper because it can be interpreted as a linearized version of two layer neural networks.

# Generalization error

**(Generalization error)** We define population risk and the empirical risk of a predictor  $\mathbf{a} \in \mathbb{R}^N$  as

$$R(\mathbf{a}) = \mathbb{E}_{\mathbf{x}, y} (y - f_{\mathbf{a}}(\mathbf{x}))^2$$
$$\widehat{R}_n(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\mathbf{a}}(\mathbf{x}_i))^2$$

and the empirical risk minimizer with vanishing regularization as

$$\mathbf{a}_{\min} = \lim_{\lambda \rightarrow 0^+} \arg \min_{\mathbf{a}} \left[ \widehat{R}_n(\mathbf{a}) + \lambda \|\mathbf{a}\|_2^2 \right].$$

In the overparameterized regime ( $N > n$ ), under mild conditions, we have  $\min_{\mathbf{a}} \widehat{R}_n(\mathbf{a}) = \widehat{R}_n(\mathbf{a}_{\min}) = 0$ . In this regime,  $\mathbf{a}_{\min}$  can be interpreted as the minimum  $\ell_2$  norm interpolator.

$\Rightarrow R(\mathbf{a}_{\min})$  : Generalization error

# Uniform convergence bounds

**(Uniform convergence bounds)** We denote the uniform convergence bound over a norm ball and the uniform convergence over interpolators in the norm ball by

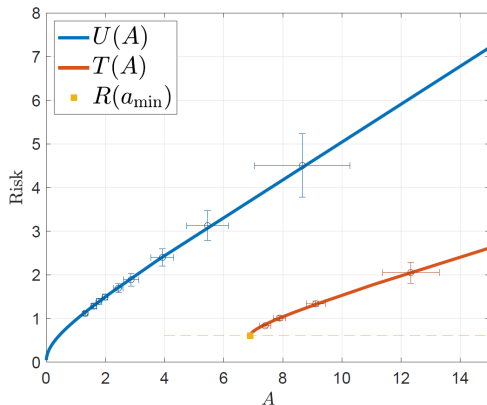
$$U(A, N, n, d) \equiv \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} \left( R(\mathbf{a}) - \widehat{R}_n(\mathbf{a}) \right)$$
$$T(A, N, n, d) \equiv \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A, \widehat{R}_n(\mathbf{a})=0} R(\mathbf{a})$$

Note that we need  $A \geq (N/d) \|\mathbf{a}_{\min}\|_2^2$  to satisfy  $\{\mathbf{a} \in \mathbb{R}^N : (N/d)\|\mathbf{a}\|_2^2 \leq A, \widehat{R}_n(\mathbf{a}) = 0\} \neq \emptyset$ .

We want to quantify  $U(A), T(A), R(\mathbf{a}_{\min})$  (abuse of notation)

# Uniform convergence bounds and Generalization error

We want to quantify  $U(A), T(A), R(\mathbf{a}_{\min})$  (abuse of notation)



$\sigma : \text{Relu}, (N, n, d) = (500, 300, 200)$

$$U(A) \geq T(A) \geq R(\mathbf{a}_{\min}) \text{ for any } A \geq (N/d) \|\mathbf{a}_{\min}\|_2^2$$

## Asymptotic limit ( $d \rightarrow \infty$ )

We approach this problem in the limit  $d \rightarrow \infty$  with  $N/d \rightarrow \psi_1$  (the number of parameters) and  $n/d \rightarrow \psi_2$  (sample size). We further assume the setting of a linear target function  $f_d$  and a nonlinear activation function  $\sigma$ .

In this regime, our main result Theorem 1 will show that,

$$\begin{aligned}U(A, N, n, d) &\xrightarrow{d \rightarrow \infty} \mathcal{U}(A, \psi_1, \psi_2), \\T(A, N, n, d) &\xrightarrow{d \rightarrow \infty} \mathcal{T}(A, \psi_1, \psi_2),\end{aligned}$$

and Theorem 1 of Mei & Montanari (2019) shows

$$\begin{aligned}(N/d) \|\mathbf{a}_{\min}\|_2^2 &\xrightarrow{d \rightarrow \infty} \mathcal{A}(\psi_1, \psi_2), \\R(\mathbf{a}_{\min}) &\xrightarrow{d \rightarrow \infty} \mathcal{R}(\psi_1, \psi_2).\end{aligned}$$



## Asymptotic limit ( $d \rightarrow \infty + \psi_1 \rightarrow \infty$ )

In the limit  $\psi_1 \rightarrow \infty$  after  $d \rightarrow \infty$ , we denote

$$\mathcal{U}_\infty(A, \psi_2) \equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{U}(A, \psi_1, \psi_2)$$

$$\mathcal{T}_\infty(A, \psi_2) \equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{T}(A, \psi_1, \psi_2)$$

$$\mathcal{A}_\infty(\psi_2) \equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{A}(\psi_1, \psi_2)$$

$$\mathcal{R}_\infty(\psi_2) \equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{R}(\psi_1, \psi_2)$$

## Asymptotic limit (substituting $A$ )

In order for  $\mathcal{U}$  and  $\mathcal{T}$  to serve as proper bounds for  $R(\mathbf{a}_{\min})$ , we need to take at least  $A \geq \psi_1 \|\mathbf{a}_{\min}\|_2^2$ . Therefore, we will choose  $A = \alpha \psi_1 \|\mathbf{a}_{\min}\|_2^2$  for some  $\alpha > 1$ .

Define

$$\begin{aligned}\mathcal{U}^{(\alpha)}(\psi_1, \psi_2) &\equiv \mathcal{U}(\alpha A(\psi_1, \psi_2), \psi_1, \psi_2), \\ \mathcal{T}^{(\alpha)}(\psi_1, \psi_2) &\equiv \mathcal{T}(\alpha A(\psi_1, \psi_2), \psi_1, \psi_2),\end{aligned}$$

and

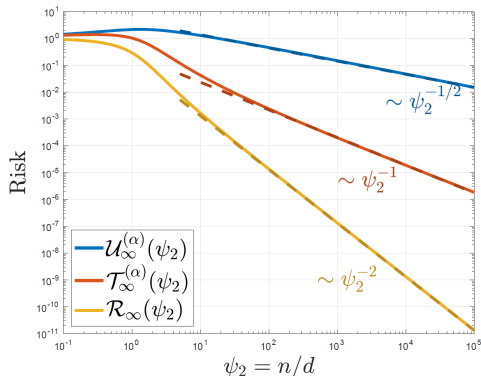
$$\begin{aligned}\mathcal{U}_{\infty}^{(\alpha)}(\psi_2) &\equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{U}^{(\alpha)}(\psi_1, \psi_2), \\ \mathcal{T}_{\infty}^{(\alpha)}(\psi_2) &\equiv \lim_{\psi_1 \rightarrow \infty} \mathcal{T}^{(\alpha)}(\psi_1, \psi_2).\end{aligned}$$

# Asymptotic power laws ( $d \rightarrow \infty + \mathcal{A}_\infty + \text{noiseless}$ )

$$\mathcal{U}_\infty^{(\alpha)}(\psi_2; \tau^2 = 0) \sim \psi_2^{-1/2}$$

$$\mathcal{T}_\infty^{(\alpha)}(\psi_2; \tau^2 = 0) \sim \psi_2^{-1}$$

$$\mathcal{R}_\infty^{(\alpha)}(\psi_2; \tau^2 = 0) \sim \psi_2^{-2}$$



## Asymptotic power laws ( $d \rightarrow \infty + \mathcal{A}_\infty + \text{noisy}$ )

$$\mathcal{U}_\infty^{(\alpha)}(\psi_2; \tau^2) - \tau^2 \sim \psi_2^{1/2}$$

$$\mathcal{T}_\infty^{(\alpha)}(\psi_2; \tau^2) - \tau^2 \sim 1$$

$$\mathcal{R}_\infty(\psi_2; \tau^2) - \tau^2 \sim \psi_2^{-1}$$

In the presence of label noise, the excess risk  $\mathcal{R}_\infty - \tau^2$  vanishes in the large sample limit. In contrast, the classical uniform convergence  $\mathcal{U}_\infty$  becomes vacuous, whereas the uniform convergence over interpolators  $\mathcal{T}_\infty$  converges to a constant, which gives a non-vacuous bound of  $\mathcal{R}_\infty$ .

## Asymptotic power laws ( $d \rightarrow \infty$ )

The inferred asymptotic law gives (c.f. Figure 4)

$$\mathcal{U}^{(\alpha)}(\psi_1, \psi_2) - \mathcal{U}_{\infty}^{(\alpha)}(\psi_2) \sim \psi_1^{-1}$$

$$\mathcal{T}^{(\alpha)}(\psi_1, \psi_2) - \mathcal{T}_{\infty}^{(\alpha)}(\psi_2) \sim \psi_1^{-1}$$

$$\mathcal{R}(\psi_1, \psi_2) - \mathcal{R}_{\infty}(\psi_2) \sim \psi_1^{-1}$$

Note that large  $\psi_1$  should be interpreted as the model being heavily overparametrized (a large width network). This asymptotic power law implies that, both uniform convergence bounds correctly predict the decay of the test error with the increase of the number of features.

# Conclusion

- They focus on investigating uniform convergence bound( $U$ ), uniform convergence bound over interpolators( $T$ ) and generalization error( $R$ ).
- The contribution of this paper is that they provide that the asymptotic limit of  $U, T$  that can be calculated through proof.
- They conjecture asymptotic power laws using simulations(figure).
- Their results exhibit a setting in which standard uniform convergence bound is vacuous while uniform convergence over interpolators gives a non-trivial bound of the actual generalization error.

- Yang, Z., Bai, Y., Mei, S. (2021). Exact gap between generalization error and uniform convergence in random feature models. arXiv preprint arXiv:2103.04554.
- Mei, S., Montanari, A. (2019). The Generalization Error of Random Features Regression: Precise Asymptotics and the Double Descent Curve. Communications on Pure and Applied Mathematics.