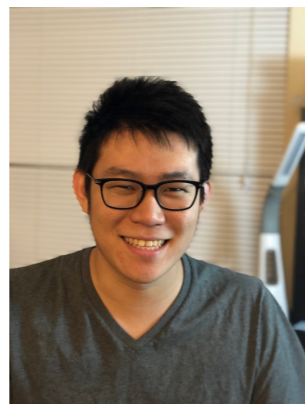




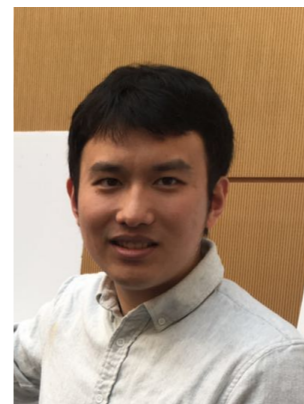
# Exact Gap between Generalization Error and Uniform Convergence\*

---

*Speaker: Zitong Yang*



*Yu Bai*



*Song Mei*

*ICML 2021*

\* <https://arxiv.org/pdf/2103.04554.pdf>

# Uniform Convergence Puzzle

- Given a training set of size  $n$

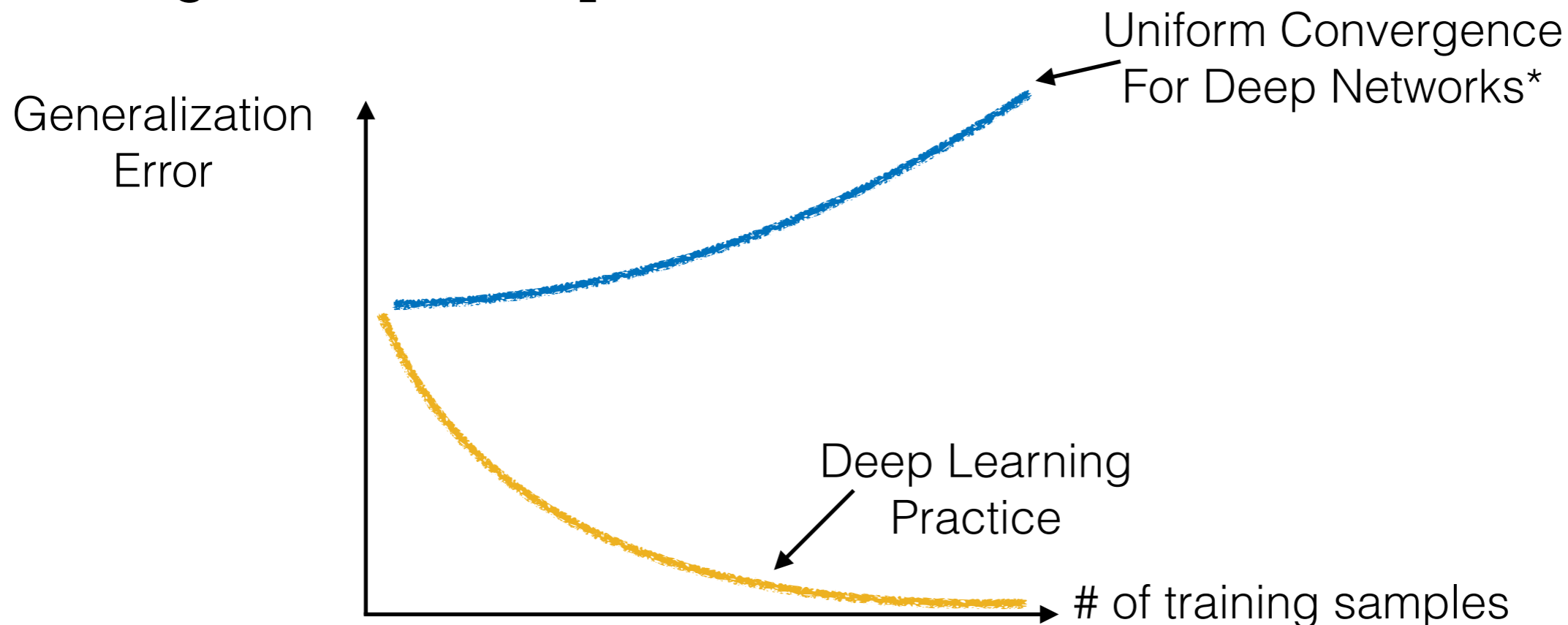
$$\underbrace{R(f) - \hat{R}_n(f)}_{\text{Generalization Error}} \leq \overbrace{\sup_{f \in \mathcal{F}} \left\{ R(f) - \hat{R}_n(f) \right\}}^{\text{Uniform Convergence}} = O\left(\sqrt{\frac{\text{Complexity of } \mathcal{F}}{n}}\right)$$

# Uniform Convergence Puzzle

- Given a training set of size  $n$

$$\underbrace{R(f) - \hat{R}_n(f)}_{\text{Generalization Error}} \leq \overbrace{\sup_{f \in \mathcal{F}} \{R(f) - \hat{R}_n(f)\}}^{\text{Uniform Convergence}} = O\left(\sqrt{\frac{\text{Complexity of } \mathcal{F}}{n}}\right)$$

- The generalization puzzle



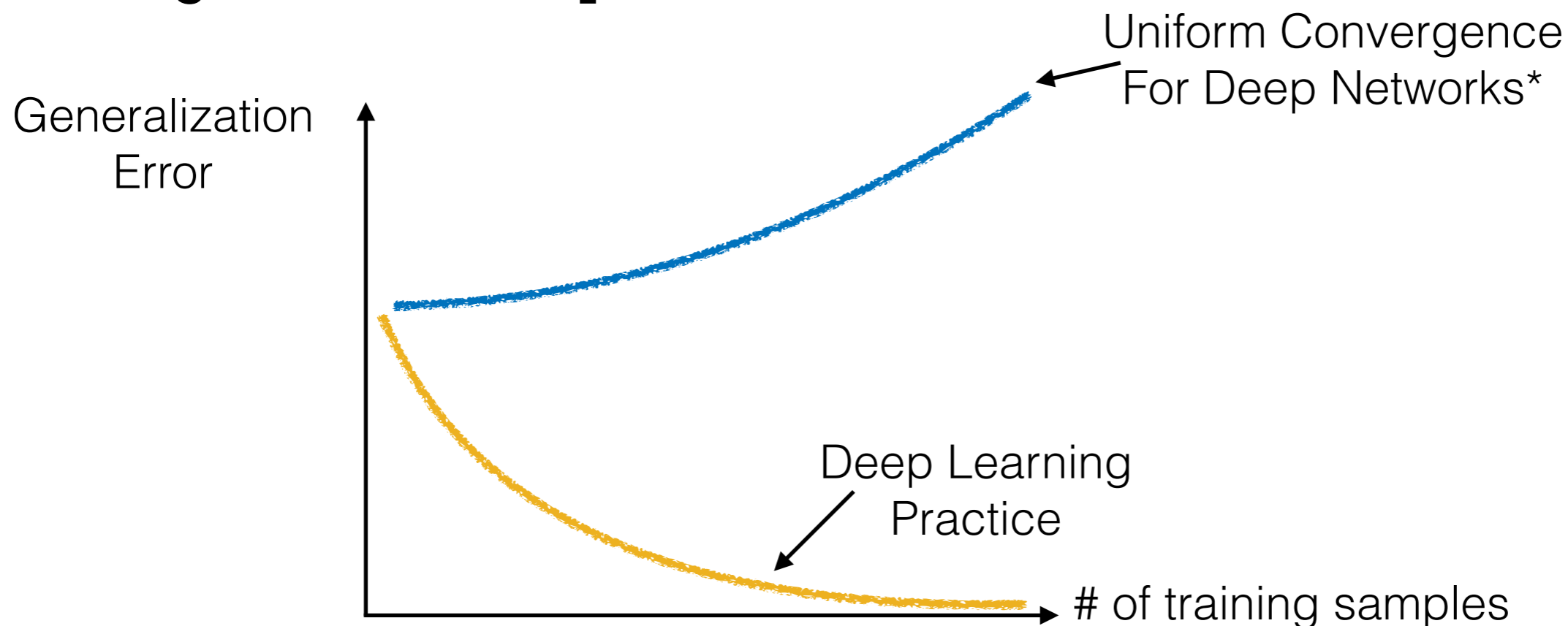
\* Nagarajan&Kolter, Uniform convergence may be unable to explain generalization

# Solution: Capacity Reduced Model Class

- Fix some  $\mathcal{F}_{\text{reduced}} \subset \mathcal{F}$

$$\underbrace{R(f) - \hat{R}_n(f)}_{\text{Generalization Error}} \leq \overbrace{\sup_{f \in \mathcal{F}_{\text{reduced}}} \{R(f) - \hat{R}_n(f)\}}^{\text{Uniform Convergence}} = O\left(\sqrt{\frac{\text{Complexity of } \mathcal{F}_{\text{reduced}}}{n}}\right)$$

- The generalization puzzle



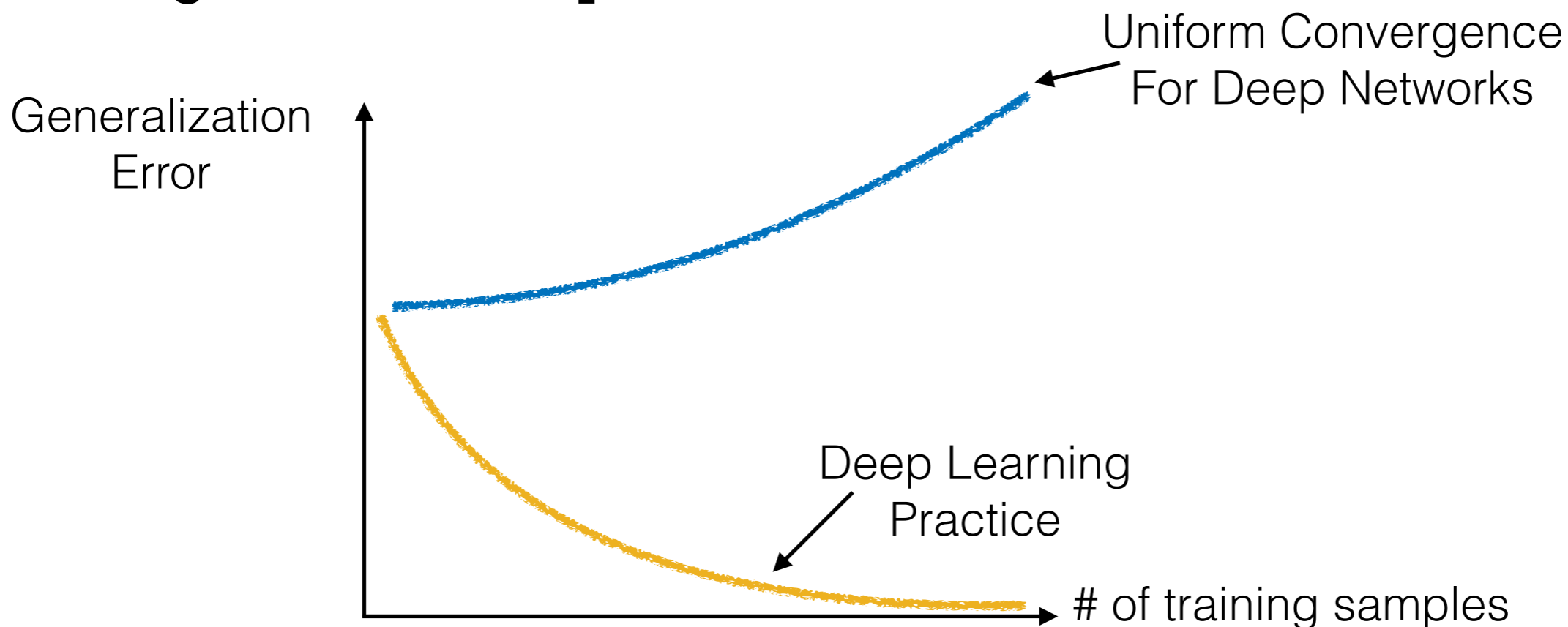
\* Nagarajan&Kolter, Uniform convergence may be unable to explain generalization

# Solution: Capacity Reduced Model Class

- As deep networks interpolates, a natural class to consider is\*

$$\mathcal{F}_{\text{reduced}} = \left\{ f \in \mathcal{F}, \hat{R}_n(f) = 0 \right\} \subset \mathcal{F}$$

- The generalization puzzle



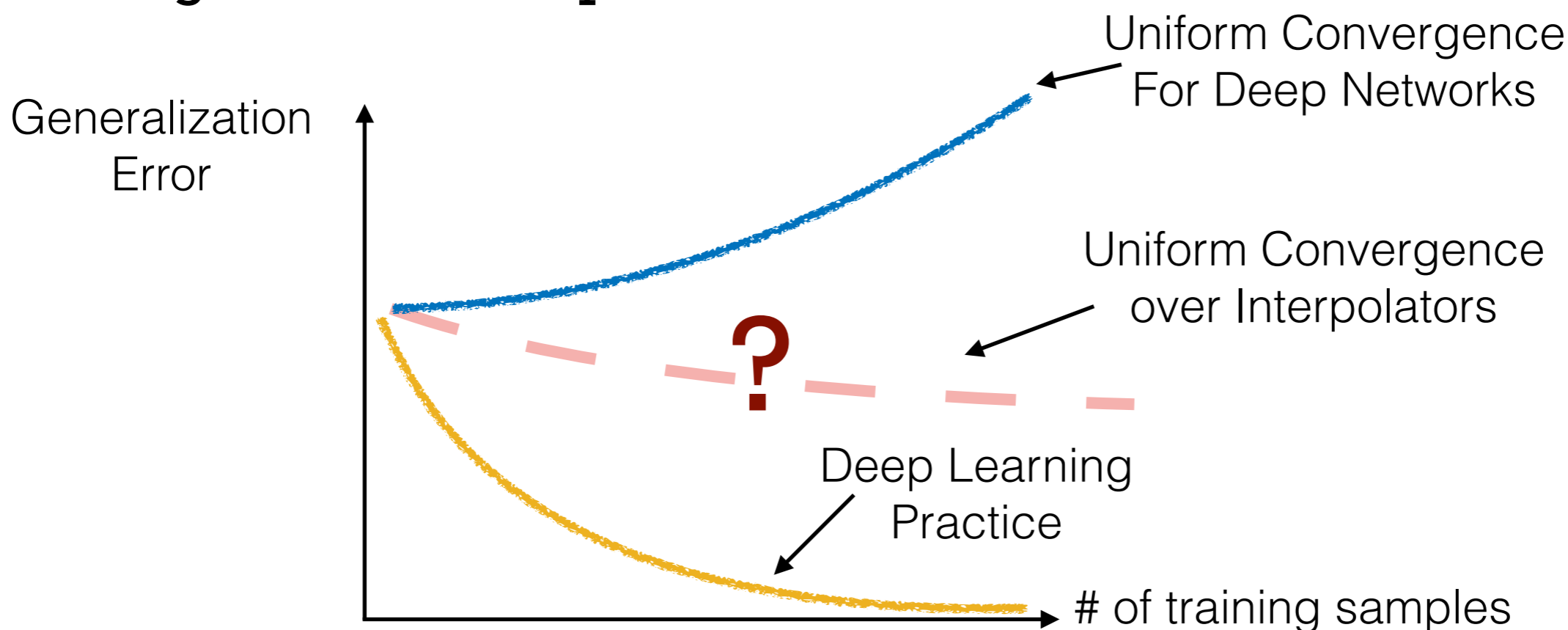
\* Zhou, Sutherland, Srebro. On uniform convergence and low-norm interpolation learning.

# Solution: Capacity Reduced Model Class

- As deep networks interpolate, a natural class to consider is\*

$$\mathcal{F}_{\text{reduced}} = \left\{ f \in \mathcal{F}, \hat{R}_n(f) = 0 \right\} \subset \mathcal{F}$$

- The generalization puzzle



\* Zhou, Sutherland, Srebro. On uniform convergence and low-norm interpolation learning.

# Solution: Capacity Reduced Model Class

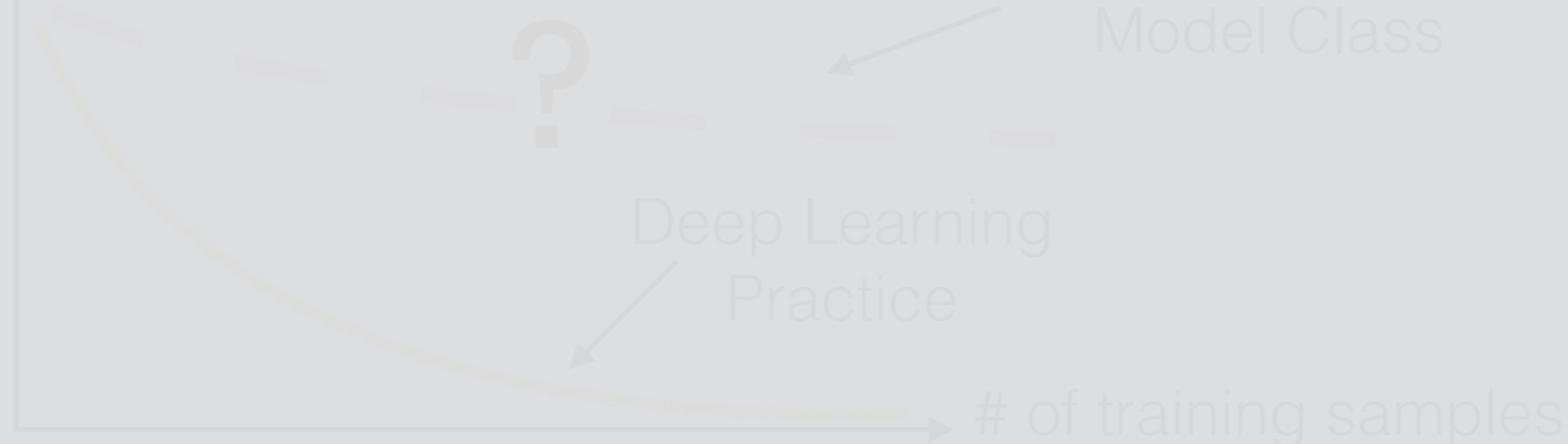
- As deep networks interpolates, a natural class to consider is\*



What's the **exact** gap between

- **Generalization Error**
- **Uniform Convergence**
- **Capacity Reduced Uniform Convergence ?**

Generalization Error



# Random Features Model

- Random Features Model: functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, A) = \left\{ f(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle), \|\mathbf{a}\|_2 \leq A, \mathbf{a} \in \mathbb{R}^N \right\},$$

where  $\boldsymbol{\theta}_j \sim \text{Unif}(\mathbb{S}^{d-1})$  are random weights.



# Random Features Model

- Random Features Model: functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, A) = \left\{ f(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle), \|\mathbf{a}\|_2 \leq A, \mathbf{a} \in \mathbb{R}^N \right\},$$

where  $\boldsymbol{\theta}_j \sim \text{Unif}(\mathbb{S}^{d-1})$  are random weights.

- Given a dataset  $(\mathbf{x}_i, y_i)_{i=1 \dots n}$ , where  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{S}^{d-1})$  and  $y_i$  are noisy observations the signal  $\boldsymbol{\theta}_\star: y_i \sim \mathcal{N}(\langle \mathbf{x}_i, \boldsymbol{\theta}_\star \rangle, \tau^2)$ , define

# Random Features Model

- Random Features Model: functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, A) = \left\{ f(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle), \|\mathbf{a}\|_2 \leq A, \mathbf{a} \in \mathbb{R}^N \right\},$$

where  $\boldsymbol{\theta}_j \sim \text{Unif}(\mathbb{S}^{d-1})$  are random weights.

- Given a dataset  $(\mathbf{x}_i, y_i)_{i=1 \dots n}$ , where  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{S}^{d-1})$  and  $y_i$  are noisy observations the signal  $\boldsymbol{\theta}_\star: y_i \sim \mathcal{N}(\langle \mathbf{x}_i, \boldsymbol{\theta}_\star \rangle, \tau^2)$ , define

$$\hat{R}_n(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \left[ y_i - \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}_i, \boldsymbol{\theta}_i \rangle) \right]^2,$$

# Random Features Model

- Random Features Model: functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, A) = \left\{ f(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle), \|\mathbf{a}\|_2 \leq A, \mathbf{a} \in \mathbb{R}^N \right\},$$

where  $\boldsymbol{\theta}_j \sim \text{Unif}(\mathbb{S}^{d-1})$  are random weights.

- Given a dataset  $(\mathbf{x}_i, y_i)_{i=1 \dots n}$ , where  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{S}^{d-1})$  and  $y_i$  are noisy observations the signal  $\boldsymbol{\theta}_\star: y_i \sim \mathcal{N}(\langle \mathbf{x}_i, \boldsymbol{\theta}_\star \rangle, \tau^2)$ , define

$$\hat{R}_n(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \left[ y_i - \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_i \rangle) \right]^2, \quad R(\mathbf{a}) = \mathbb{E}_{\mathbf{x}, y} \left[ y_i - \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_i \rangle) \right]^2.$$

# Random Features Model

- Random Features Model: functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\mathcal{F}_{\text{RF}}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N, A) = \left\{ f(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle), \|\mathbf{a}\|_2 \leq A, \mathbf{a} \in \mathbb{R}^N \right\},$$

where  $\boldsymbol{\theta}_j \sim \text{Unif}(\mathbb{S}^{d-1})$  are random weights.

- Given a dataset  $(\mathbf{x}_i, y_i)_{i=1 \dots n}$ , where  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \text{Unif}(\mathbb{S}^{d-1})$  and  $y_i$  are noisy observations the signal  $\boldsymbol{\theta}_\star: y_i \sim \mathcal{N}(\langle \mathbf{x}_i, \boldsymbol{\theta}_\star \rangle, \tau^2)$ , define

$$\hat{R}_n(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \left[ y_i - \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_i \rangle) \right]^2, \quad R(\mathbf{a}) = \mathbb{E}_{\mathbf{x}, y} \left[ y_i - \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}, \boldsymbol{\theta}_i \rangle) \right]^2.$$

- When  $N > n$ , the min-norm interpolator exists w.h.p.

$$\mathbf{a}_{\min} = \arg \min_{\hat{R}_n(\mathbf{a})=0} \|\mathbf{a}\|.$$

# Random Features Model

- Random Features Model: functions  $f: \mathbb{R}^d \rightarrow \mathbb{R}$



What's the **exact** gap between

- **Generalization Error**

- **Uniform Convergence**

- **Capacity Reduced Uniform Convergence**



- When  $N > n$ , the min-norm interpolator exists w.h.p.

$$a_{\min} = \arg \min_{\hat{R}_n(a)=0} \|a\|.$$

# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

- Uniform Convergence

$$U(A, N, n, d) = \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} R(\mathbf{a}) - \hat{R}_n(\mathbf{a})$$

# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

- Uniform Convergence

$$U(A, N, n, d) = \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} R(\mathbf{a}) - \hat{R}_n(\mathbf{a})$$

- Uniform Convergence over Interpolators

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



# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

- Uniform Convergence

$$U(A, N, n, d) = \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} R(\mathbf{a}) - \hat{R}_n(\mathbf{a})$$

- Uniform Convergence over Interpolators

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

- With the choice of  $A = \alpha \cdot (N/d)\|\mathbf{a}_{\min}\|_2^2$  (e.g.  $\alpha = 1.1$ ).

# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

- Uniform Convergence

$$U(A, N, n, d) = \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} R(\mathbf{a}) - \hat{R}_n(\mathbf{a})$$

- Uniform Convergence over Interpolators

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

- With the choice of  $A = \alpha \cdot (N/d)\|\mathbf{a}_{\min}\|_2^2$  (e.g.  $\alpha = 1.1$ ).

- Consider  $\lim_{d \rightarrow \infty} N/d = \psi_1$ ,  $\lim_{d \rightarrow \infty} n/d = \psi_2$ .

# Uniform Convergence in Random Features Model

- Generalization Error

$$R(N, n, d) = R(\mathbf{a}_{\min}) - \underbrace{\hat{R}_n(\mathbf{a}_{\min})}_{=0}$$

- Uniform Convergence

$$U(A, N, n, d) = \sup_{(N/d)\|\mathbf{a}\|_2^2 \leq A} R(\mathbf{a}) - \hat{R}_n(\mathbf{a})$$

- Uniform Convergence over Interpolators

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

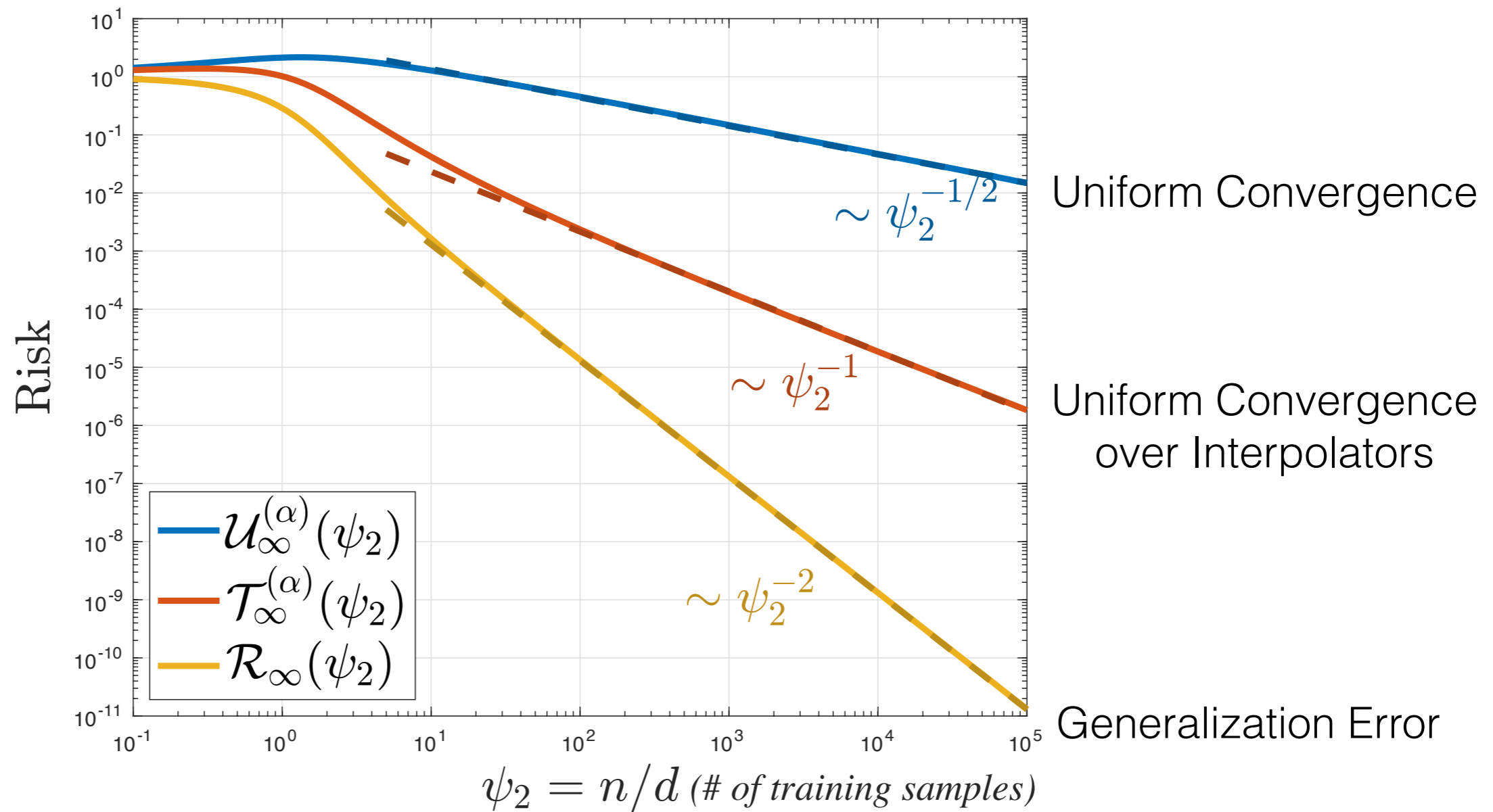
- With the choice of  $A = \alpha \cdot (N/d)\|\mathbf{a}_{\min}\|_2^2$  (e.g.  $\alpha = 1.1$ ).

- Consider  $\lim_{d \rightarrow \infty} N/d = \psi_1$ ,  $\lim_{d \rightarrow \infty} n/d = \psi_2$ .

Allows exact computation instead non-asymptotic bounds.

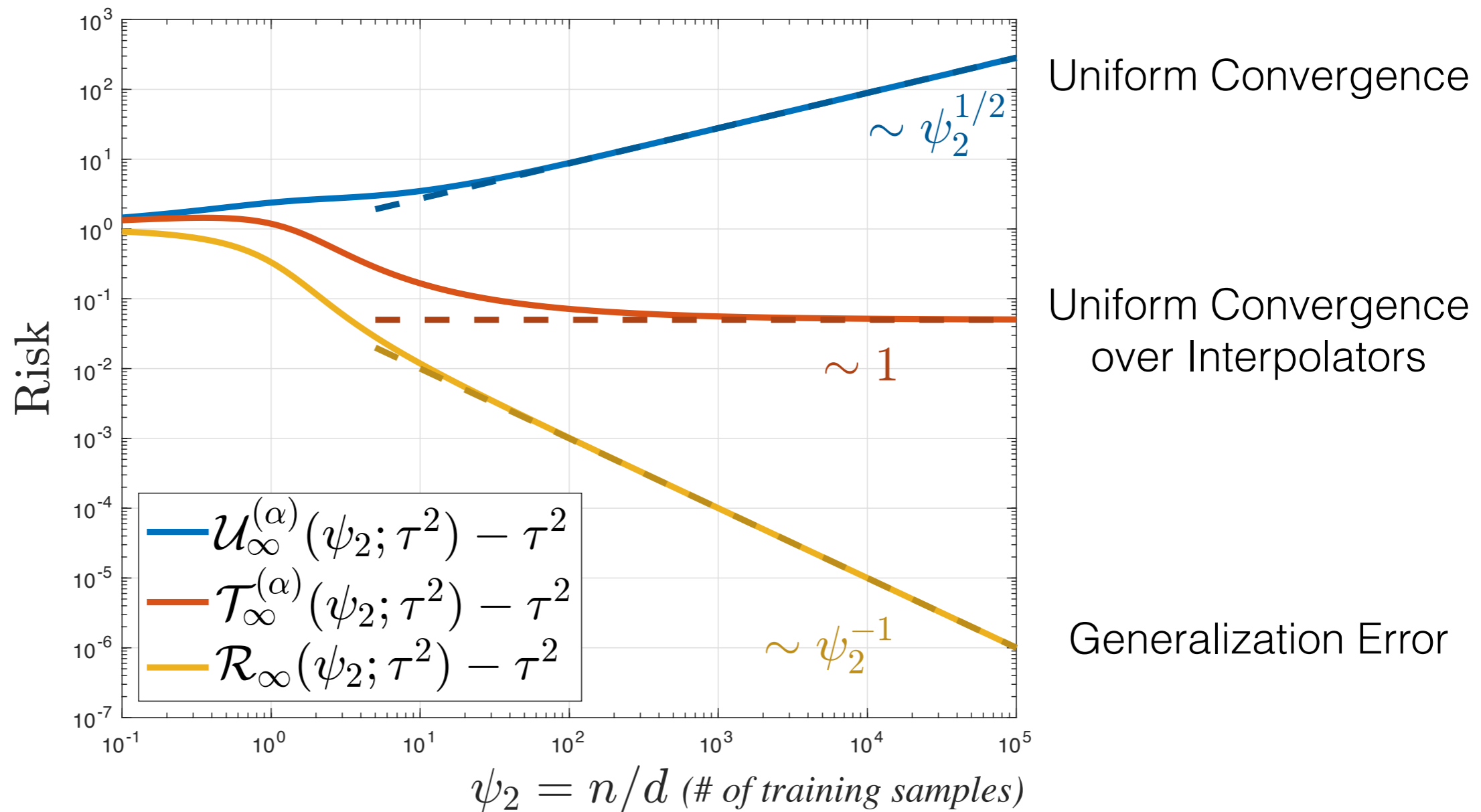
# Main Results: Exact Analytical Formulae

- In the noiseless regime  $\tau^2 = 0$



# Main Results: Exact Analytical Formulae

- In the noisy regime  $\tau^2 > 0$



**Thanks!**