

Be aware of model capacity when talking about generalization in modern machine learning

Fanghui Liu

fanghui.liu@warwick.ac.uk

Department of Computer Science, University of Warwick, UK

Centre for Discrete Mathematics and its Applications (DIMAP), Warwick

at Department of Statistics, University of Oxford



WARWICK
THE UNIVERSITY OF WARWICK



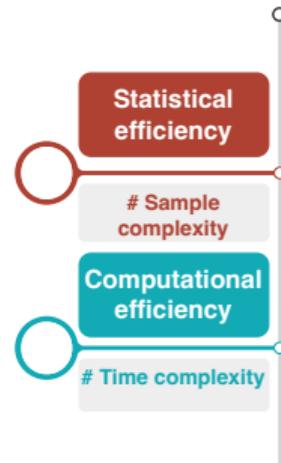
The
Alan Turing
Institute



My research

❑ Research interests

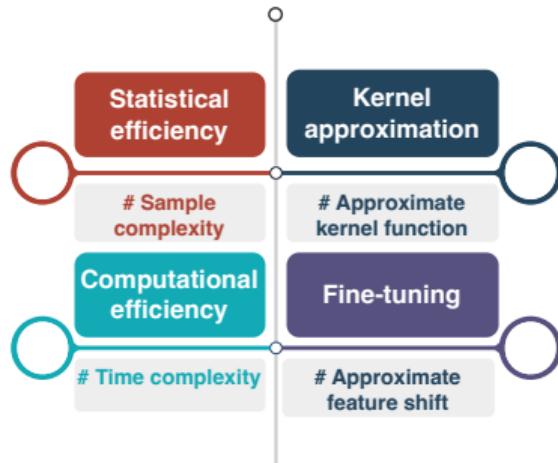
- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML



My research

❑ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML



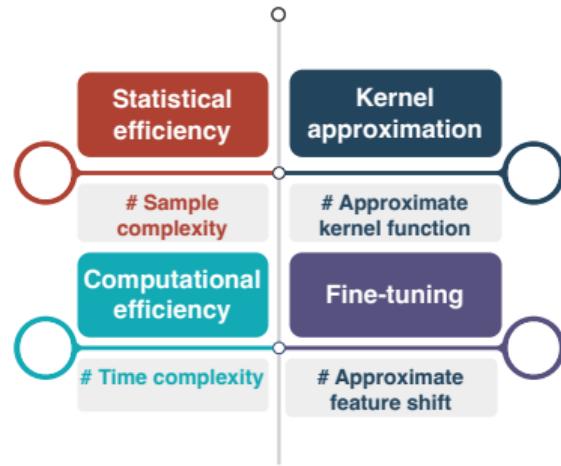
My research

❑ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

❑ Research goal

- characterize **learning efficiency** in theory
- contribute to practice



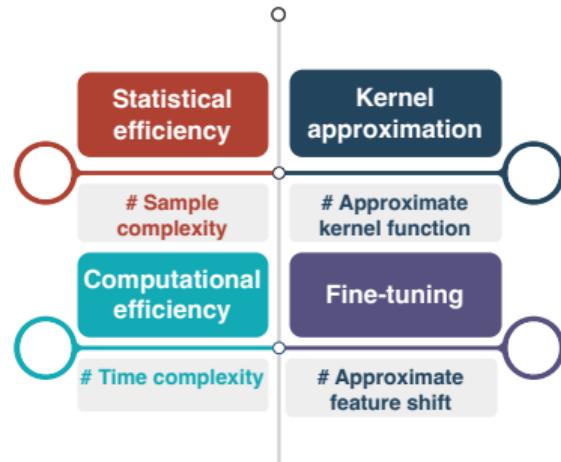
My research

Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

Research goal

- characterize **learning efficiency** in theory
- contribute to practice



Learning efficiency (Curse of Dimensionality, CoD)

Machine learning works in **high dimensions** that can be a **curse!**

— David Donoho, 2000. (Richard E. Bellman, 1957)

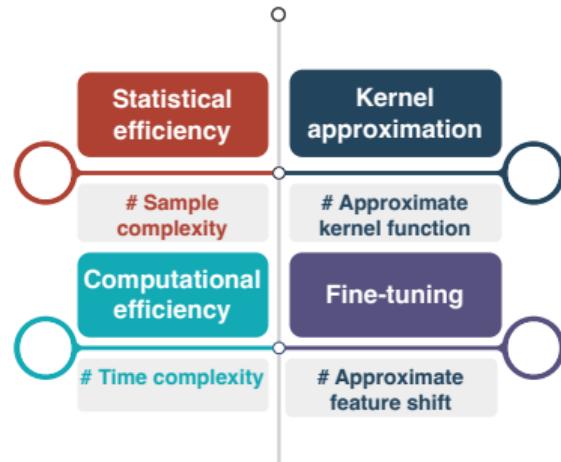
My research

Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

Research goal

- characterize **learning efficiency** in theory
- contribute to practice



Learning efficiency (Curse of Dimensionality, CoD)

Machine learning works in **high dimensions** that can be a **curse!**

— David Donoho, 2000. (Richard E. Bellman, 1957)



Data



Model



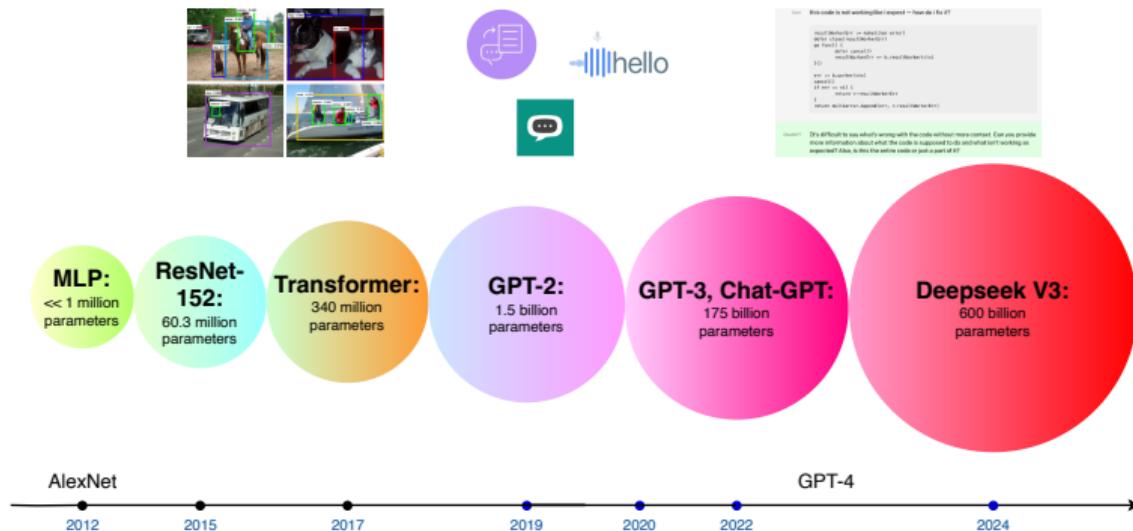
Algorithm



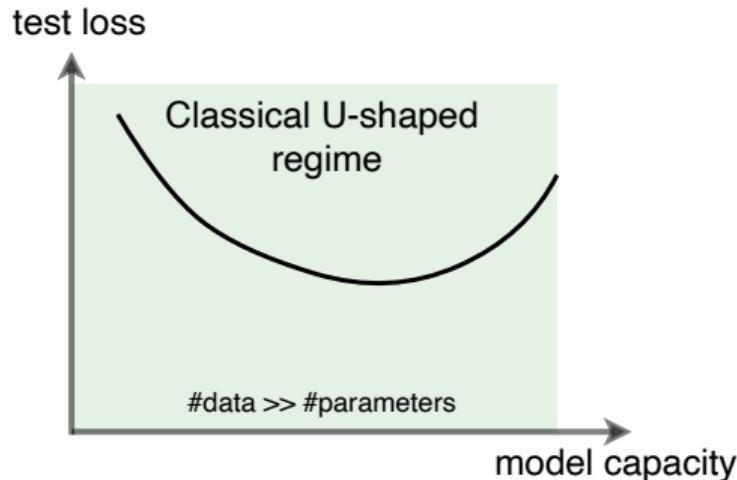
Compute

In the era of machine learning

Prefer more data and larger model to obtain better performance...

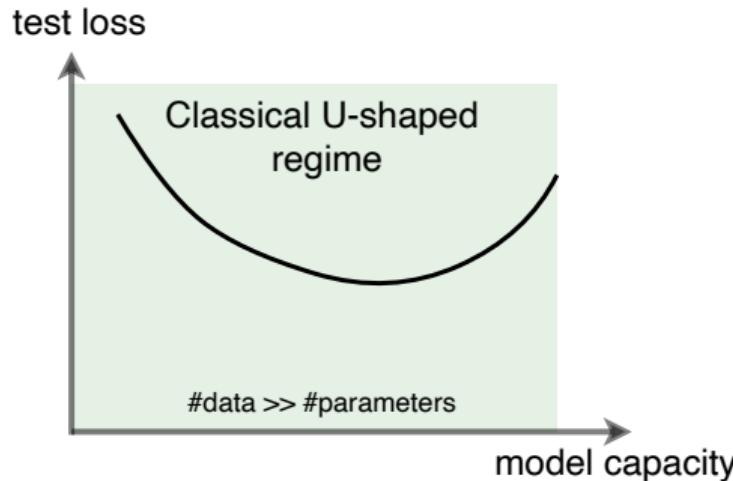


ML textbooks: Larger models tend to overfit!

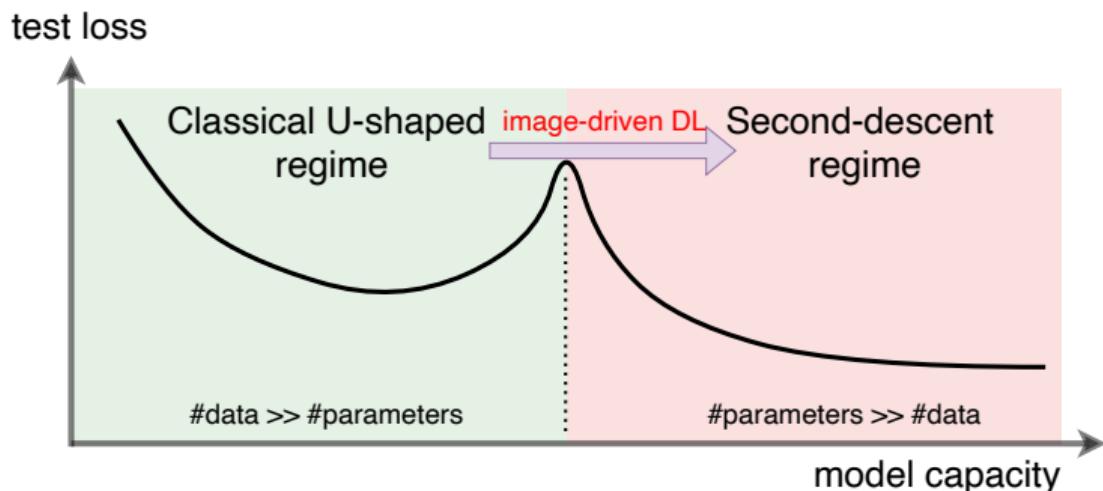


ML textbooks: Larger models tend to overfit!

Practice of deep learning: bigger models perform better!



Practice of deep learning: bigger models perform better!



Proposed explanation: double descent (Belkin et al., 2019)

Learning paradigm in the past twenty years

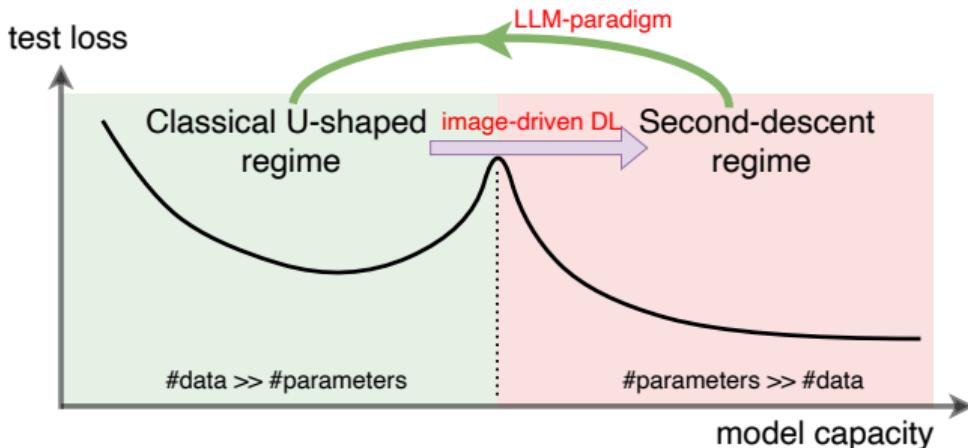


Figure 1: Paradigm among test loss, data, and model capacity.

Scaling law (Kaplan et al., 2020) in the era of LLMs

$$\text{test loss} = A \times \text{Model Size}^{-a} + B \times \text{Data Size}^{-b} + C$$

Learning paradigm in the past twenty years

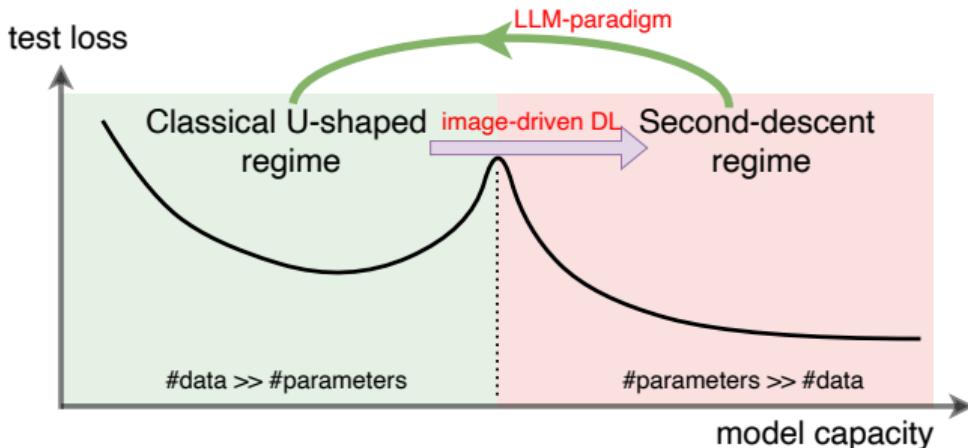


Figure 1: Paradigm among test loss, data, and model capacity.

Scaling law (Kaplan et al., 2020) in the era of LLMs

$$\text{test loss} = A \times \text{Model Size}^{-a} + B \times \text{Data Size}^{-b} + C$$

A fundamental concept in machine learning: model capacity

Too many learning curves...

- U-shaped curve (bias-variance trade-offs) ([Vapnik, 1995; Hastie et al., 2009](#))
- double (multiple) descent ([Belkin et al., 2019; Liang et al., 2020](#))
- scaling law ([Kaplan et al., 2020; Paquette et al., 2024](#))

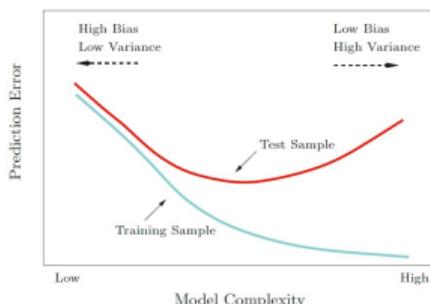
A fundamental concept in machine learning: model capacity

Too many learning curves...

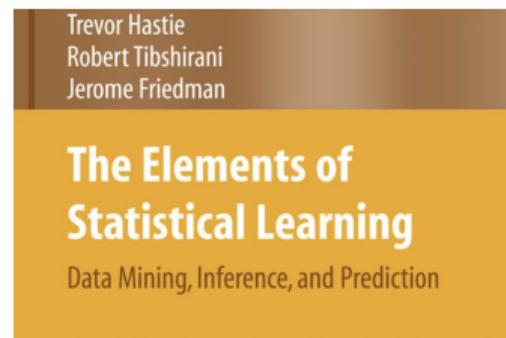
- U-shaped curve (bias-variance trade-offs) ([Vapnik, 1995; Hastie et al., 2009](#))
- double (multiple) descent ([Belkin et al., 2019; Liang et al., 2020](#))
- scaling law ([Kaplan et al., 2020; Paquette et al., 2024](#))

Bias-variance decomposition

$$\text{Test error} = \text{Bias}^2 + \text{Variance}$$



([Hastie et al., 2009](#), Figure 2.11)



A fundamental concept in machine learning: model capacity

Too many learning curves...

- U-shaped curve (bias-variance trade-offs) (Vapnik, 1995; Hastie et al., 2009)
- double (multiple) descent (Belkin et al., 2019; Liang et al., 2020)
- scaling law (Kaplan et al., 2020; Paquette et al., 2024)

Bias-variance decomposition

$$\text{Test error} = \text{Bias}^2 + \text{Variance}$$

"Remove bias-variance trade-offs from ML textbooks"

Trade-off is a **misnomer**, by Geman et al. (1992); Neal (2019); Wilson (2025).

I can define **model capacity** at random and see whatever curve I want to see.

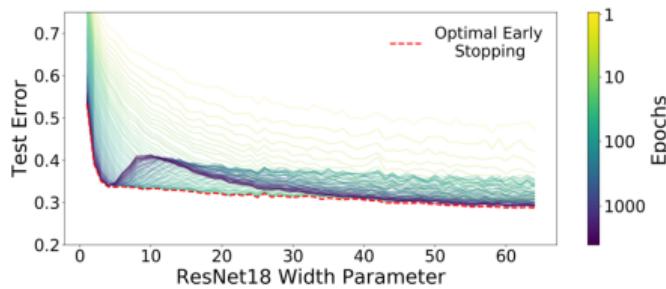
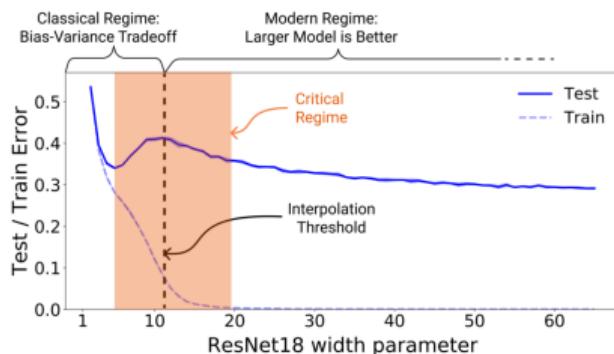
— Ben Recht, 2025

A fundamental concept in machine learning: model capacity

Too many learning curves...

- U-shaped curve (bias-variance trade-offs) (Vapnik, 1995; Hastie et al., 2009)
- double (multiple) descent (Belkin et al., 2019; Liang et al., 2020)
- scaling law (Kaplan et al., 2020; Paquette et al., 2024)

Double descent can disappear for the same architecture!



(a) Results on ResNet18 (Nakkiran et al., 2019) (b) Optimal early stopping (Nakkiran et al., 2019).

Today's talk: Learning with norm-based capacity

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- Theoretical studies (Neyshabur et al., 2015; Savarese et al., 2019)
- Min-norm solution (Hastie et al., 2022)
- Applications: neural networks pruning (Molchanov et al., 2017), lottery ticket hypothesis (Frankle and Carbin, 2019)

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- Theoretical studies (Neyshabur et al., 2015; Savarese et al., 2019)
- Min-norm solution (Hastie et al., 2022)
- Applications: neural networks pruning (Molchanov et al., 2017), lottery ticket hypothesis (Frankle and Carbin, 2019)

How these learning curves behave under a more suitable model capacity?

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- Theoretical studies (Neyshabur et al., 2015; Savarese et al., 2019)
- Min-norm solution (Hastie et al., 2022)
- Applications: neural networks pruning (Molchanov et al., 2017), lottery ticket hypothesis (Frankle and Carbin, 2019)

How these learning curves behave under a more suitable model capacity?

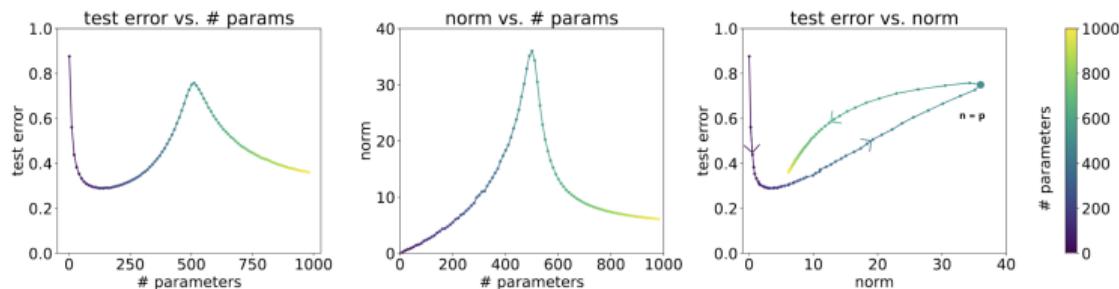


Figure 3: Stanford CS229 lecture notes (Ng and Ma, 2023, Figure 8.12).

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- How to precisely characterize the relationship under norm-based model capacity?
 - Reshape bias-variance trade-offs, double descent, scaling law under ℓ_2 norm-based capacity!
 - Yichen Wang, Yudong Chen, Lorenzo Rosasco, Fanghui Liu. *The shape of generalization through the lens of norm-based capacity control.* 2025. [arXiv](#)

Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- How to precisely characterize the relationship under norm-based model capacity?
 - Reshape bias-variance trade-offs, double descent, scaling law under ℓ_2 norm-based capacity!
 - Yichen Wang, Yudong Chen, Lorenzo Rosasco, Fanghui Liu. *The shape of generalization through the lens of norm-based capacity control*. 2025. [arXiv](#)
- What is the induced function space and statistical/computational efficiency under norm-based capacity?

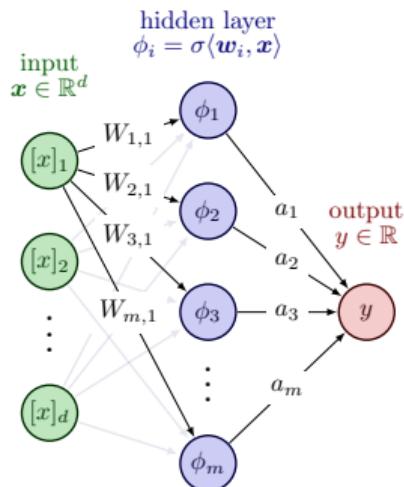
Today's talk: Learning with norm-based capacity

(Bartlett, 1998)

"The size of the weights is more important than the size of the network!"

- How to precisely characterize the relationship under norm-based model capacity?
 - Reshape bias-variance trade-offs, double descent, scaling law under ℓ_2 norm-based capacity!
 - Yichen Wang, Yudong Chen, Lorenzo Rosasco, Fanghui Liu. *The shape of generalization through the lens of norm-based capacity control*. 2025. [arXiv](#)
- What is the induced function space and statistical/computational efficiency under norm-based capacity?
 - Which function class can be efficiently learned by neural networks?
 - Fanghui Liu, Leello Dadi, and Volkan Cevher. *Learning with norm constrained, over-parameterised, two-layer neural networks*. JMLR 2024.

Background: Random features ridge regression



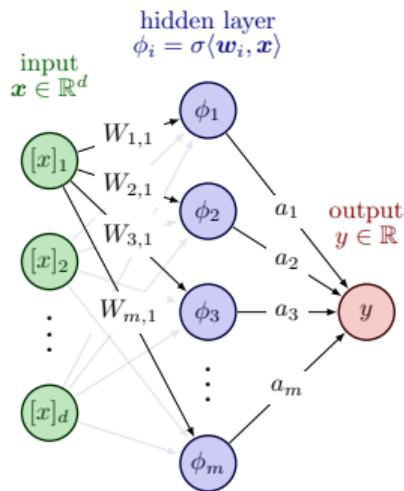
$$f_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^m \textcolor{red}{a}_i \phi(\mathbf{x}, \mathbf{w}_i), \quad \boldsymbol{\theta} := \{(a_i, \mathbf{w}_i)\}_{i=1}^m$$

- $\phi : \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$, e.g., ReLU:
 $\phi(\mathbf{x}, \mathbf{w}) = \max(\langle \mathbf{x}, \mathbf{w} \rangle, 0)$
- Random features models (RFMs) Rahimi and Recht (2007):
 - $\{\mathbf{w}_i\}_{i=1}^m \stackrel{iid}{\sim} \mu$ for a given $\mu \in \mathcal{P}(\mathcal{W})$
 - only train the second layer

$$\hat{\mathbf{a}} := \underset{\mathbf{a} \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{a}))^2 + \lambda \|\mathbf{a}\|_2^2 \right\} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{y}.$$

- $\mathbf{Z} \in \mathbb{R}^{n \times p}$ with $[\mathbf{Z}]_{ij} = \frac{1}{\sqrt{p}} \phi(\mathbf{x}_i, \mathbf{w}_j)$.
- Norm over the first-layer (untrained) $\|\mathbf{W}\|_{\text{F}}$
- Norm over the second-layer $\|\hat{\mathbf{a}}\|_2^2$

Background: Random features ridge regression



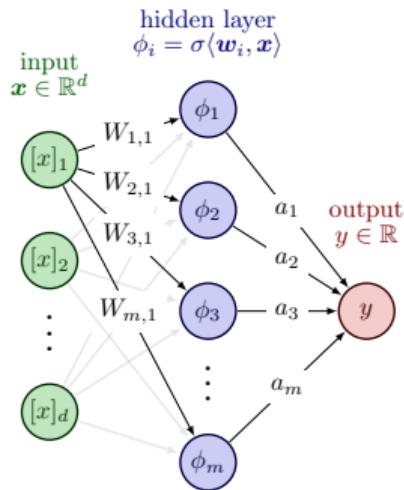
$$f_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^m a_i \phi(\mathbf{x}, \mathbf{w}_i), \quad \boldsymbol{\theta} := \{(a_i, \mathbf{w}_i)\}_{i=1}^m$$

- $\phi : \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$, e.g., ReLU:
$$\phi(\mathbf{x}, \mathbf{w}) = \max(\langle \mathbf{x}, \mathbf{w} \rangle, 0)$$
- Random features models (RFMs) Rahimi and Recht (2007):
 - $\{\mathbf{w}_i\}_{i=1}^m \stackrel{iid}{\sim} \mu$ for a given $\mu \in \mathcal{P}(\mathcal{W})$
 - only train the second layer

$$\hat{\mathbf{a}} := \operatorname{argmin}_{\mathbf{a} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{a}))^2 + \lambda \|\mathbf{a}\|_2^2 \right\} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{y}.$$

- $\mathbf{Z} \in \mathbb{R}^{n \times p}$ with $[\mathbf{Z}]_{ij} = \frac{1}{\sqrt{p}} \phi(\mathbf{x}_i, \mathbf{w}_j)$.
- Norm over the first-layer (untrained) $\|\mathbf{W}\|_{\text{F}}$
- Norm over the second-layer $\|\hat{\mathbf{a}}\|_2^2$

Background: Random features ridge regression



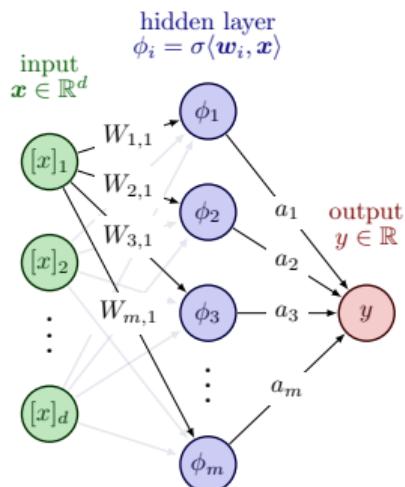
$$f_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^m \textcolor{red}{a}_i \phi(\mathbf{x}, \mathbf{w}_i), \quad \boldsymbol{\theta} := \{(a_i, \mathbf{w}_i)\}_{i=1}^m$$

- $\phi : \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$, e.g., ReLU:
$$\phi(\mathbf{x}, \mathbf{w}) = \max(\langle \mathbf{x}, \mathbf{w} \rangle, 0)$$
- Random features models (RFMs) Rahimi and Recht (2007):
 - $\{\mathbf{w}_i\}_{i=1}^m \stackrel{iid}{\sim} \mu$ for a given $\mu \in \mathcal{P}(\mathcal{W})$
 - only train the second layer

$$\hat{\mathbf{a}} := \underset{\mathbf{a} \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{a}))^2 + \lambda \|\mathbf{a}\|_2^2 \right\} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{y}.$$

- $\mathbf{Z} \in \mathbb{R}^{n \times p}$ with $[\mathbf{Z}]_{ij} = \frac{1}{\sqrt{p}} \phi(\mathbf{x}_i, \mathbf{w}_j)$.
- Norm over the first-layer (untrained) $\|\mathbf{W}\|_{\text{F}}$
- Norm over the second-layer $\|\hat{\mathbf{a}}\|_2^2$

Background: Random features ridge regression

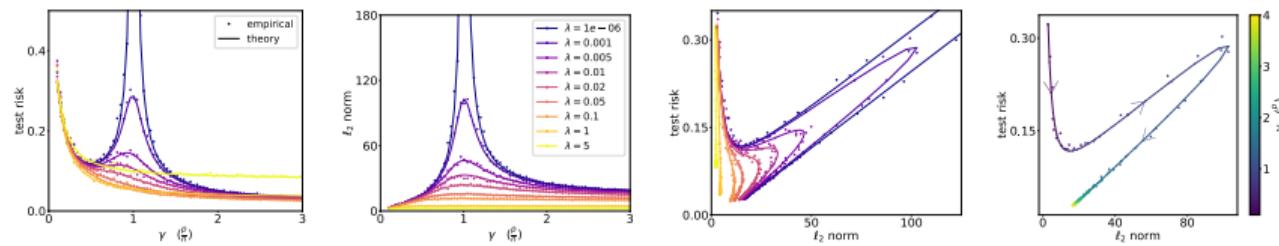
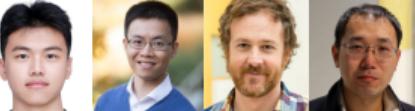


$$f_m(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^m a_i \phi(\mathbf{x}, \mathbf{w}_i), \quad \boldsymbol{\theta} := \{(a_i, \mathbf{w}_i)\}_{i=1}^m$$

- $\phi : \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$, e.g., ReLU:
$$\phi(\mathbf{x}, \mathbf{w}) = \max(\langle \mathbf{x}, \mathbf{w} \rangle, 0)$$
- Random features models (RFMs) Rahimi and Recht (2007):
 - $\{\mathbf{w}_i\}_{i=1}^m \stackrel{iid}{\sim} \mu$ for a given $\mu \in \mathcal{P}(\mathcal{W})$
 - only train the second layer

$$\hat{\mathbf{a}} := \underset{\mathbf{a} \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{a}))^2 + \lambda \|\mathbf{a}\|_2^2 \right\} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{y}.$$

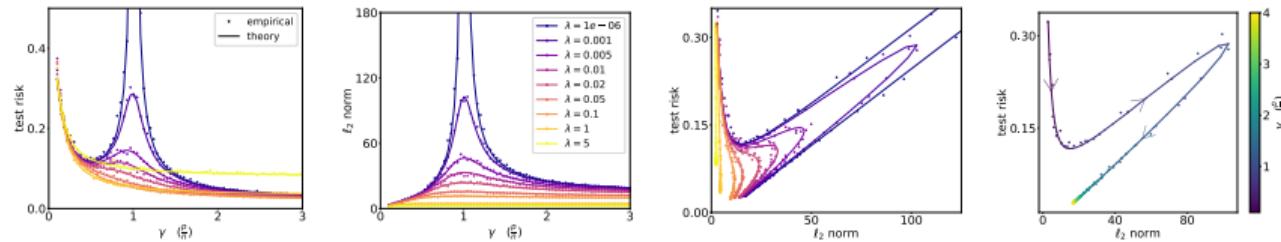
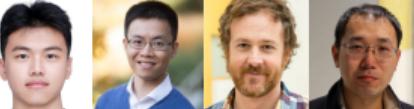
- $\mathbf{Z} \in \mathbb{R}^{n \times p}$ with $[\mathbf{Z}]_{ij} = \frac{1}{\sqrt{p}} \phi(\mathbf{x}_i, \mathbf{w}_j)$.
- Norm over the first-layer (untrained) $\|\mathbf{W}\|_{\text{F}}$
- Norm over the second-layer $\|\hat{\mathbf{a}}\|_2^2$

(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ

(c) Test Risk vs. norm

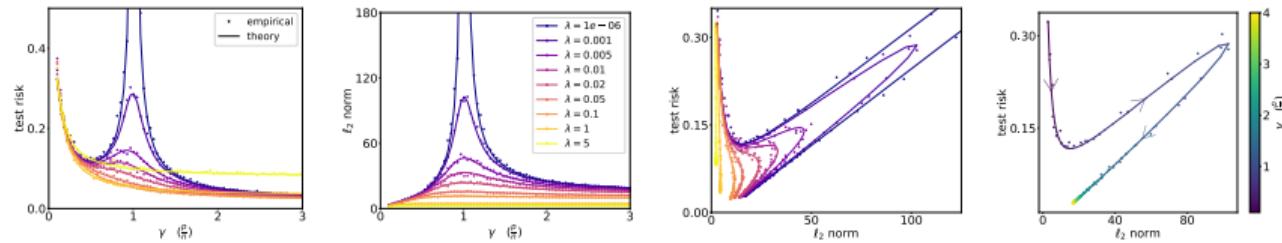
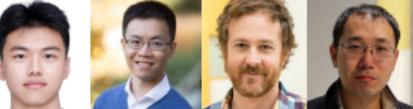
(d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size



(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ (c) Test Risk vs. norm (d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size
- Phase transition exists but double descent does not exist
- More close to **U-shaped** instead of double descent

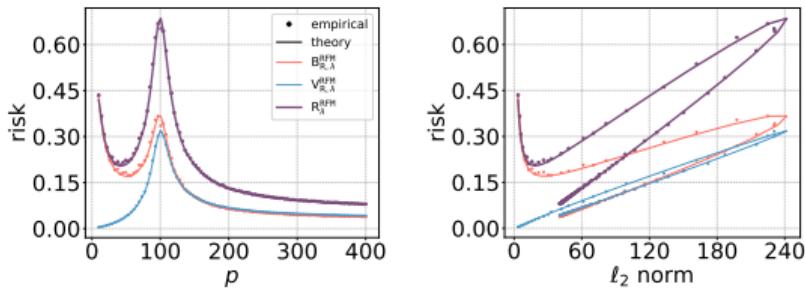
(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ

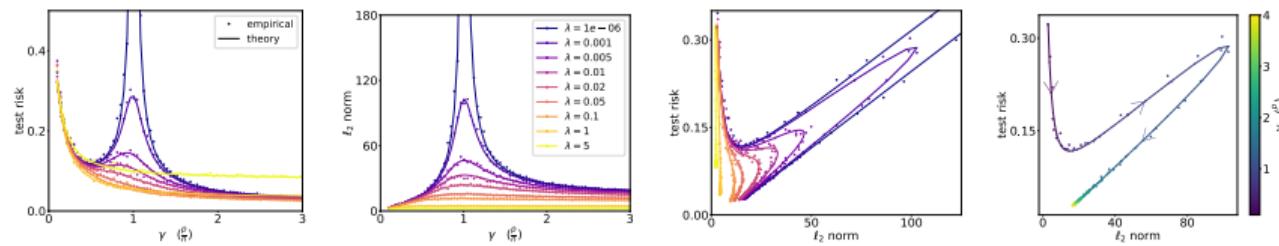
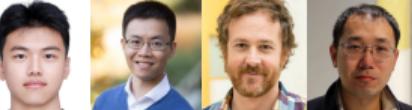
(c) Test Risk vs. norm

(d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size

$$\text{Test error} = \text{Bias}^2 + \text{Variance}$$

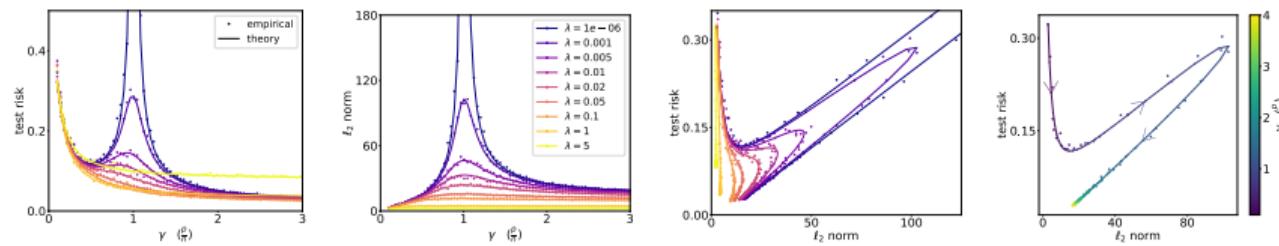
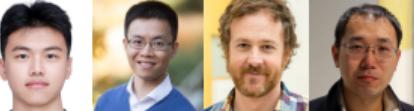


(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ

(c) Test Risk vs. norm

(d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size
- Over-parameterization is still **better than** under-parameterization

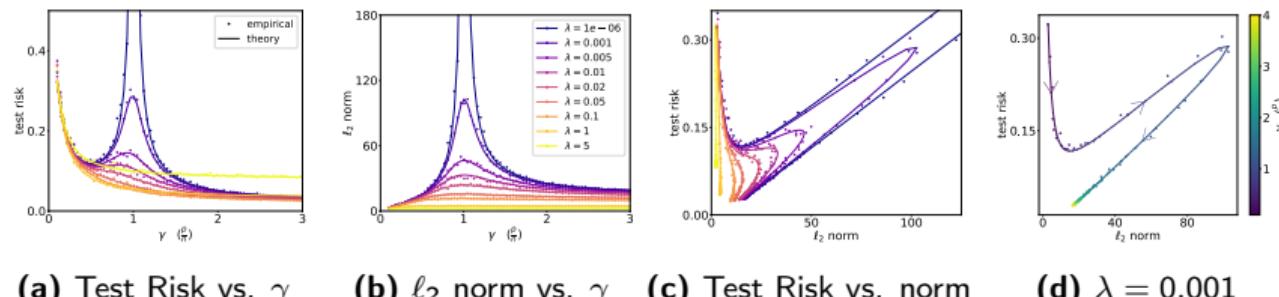
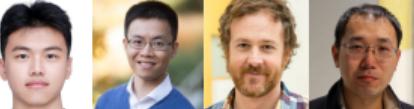
(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ

(c) Test Risk vs. norm

(d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size
- Over-parameterization is still **better than** under-parameterization
- Reshape scaling-law:

$$\text{test loss} = A \times \text{Data Size}^{-a} + B \times \text{Model Size}^{-b} + C$$
 with $a, b > 0$



(a) Test Risk vs. γ (b) ℓ_2 norm vs. γ (c) Test Risk vs. norm (d) $\lambda = 0.001$

- $\gamma := p/n$, p : model size (width), n : data size
- Over-parameterization is still **better than** under-parameterization
- Reshape scaling-law:
 $\text{test loss} = A \times \text{Data Size}^{-a} + B \times \text{Model Size}^{-b} + C$ with $a, b > 0$
 $\text{test loss} = A \times \text{Data Size}^{-a} \times \text{Norm Capacity}^{-b}$ with $a > 0$ and $b \in \mathbb{R}$

Control norm by tuning λ : L-curve (Hansen, 1992)

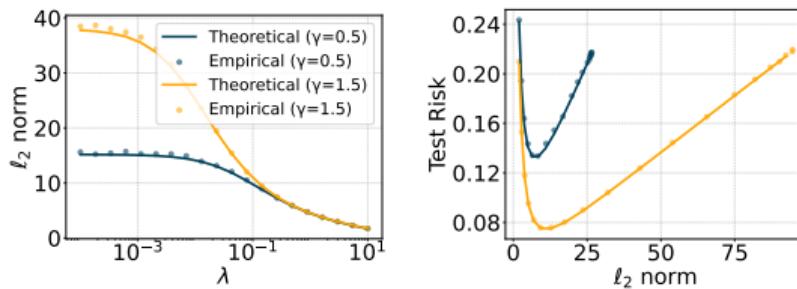
Explicit (model size) vs. Implicit (norm)

One-to-one mapping between norm and λ

Control norm by tuning λ : L-curve (Hansen, 1992)

Explicit (model size) vs. Implicit (norm)

One-to-one mapping between norm and λ



(a) Norm vs. λ (varying λ)

(b) Risk vs. Norm (varying λ)

An example of linear regression: Textbook level and beyond

- n i.i.d. samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
- $y = \langle \boldsymbol{\beta}_*, \mathbf{x} \rangle + \varepsilon$, $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{V}(\varepsilon) = \sigma^2$, covariance matrix $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$
- ridge regression: $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Target: precise analysis

The expected test risk $\mathbb{E}_\varepsilon \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\|_\Sigma^2$ vs. the norm $\mathbb{E}_\varepsilon \|\hat{\boldsymbol{\beta}}\|_2^2$

An example of linear regression: Textbook level and beyond

- n i.i.d. samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
- $y = \langle \boldsymbol{\beta}_*, \mathbf{x} \rangle + \varepsilon$, $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{V}(\varepsilon) = \sigma^2$, covariance matrix $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$
- ridge regression: $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Target: precise analysis

The expected test risk $\mathbb{E}_\varepsilon \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\|_\Sigma^2$ vs. the norm $\mathbb{E}_\varepsilon \|\hat{\boldsymbol{\beta}}\|_2^2$

An example of linear regression: Textbook level and beyond

- n i.i.d. samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
- $y = \langle \beta_*, \mathbf{x} \rangle + \varepsilon$, $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{V}(\varepsilon) = \sigma^2$, covariance matrix $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$
- ridge regression: $\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Target: precise analysis

The expected test risk $\mathbb{E}_\varepsilon \|\hat{\beta} - \beta_*\|_\Sigma^2$ vs. the norm $\mathbb{E}_\varepsilon \|\hat{\beta}\|_2^2$

- Deterministic equivalence ([Cheng and Montanari, 2024](#); [Misiakiewicz and Saeed, 2024](#)): law of large samples/dimensions in random matrix theory

The empirical spectral measure converges to a deterministic limit.

An example of linear regression: Textbook level and beyond

- n i.i.d. samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
- $y = \langle \beta_*, \mathbf{x} \rangle + \varepsilon$, $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{V}(\varepsilon) = \sigma^2$, covariance matrix $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$
- ridge regression: $\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Target: precise analysis

The expected test risk $\mathbb{E}_\varepsilon \|\hat{\beta} - \beta_*\|_\Sigma^2$ vs. the norm $\mathbb{E}_\varepsilon \|\hat{\beta}\|_2^2$

- Deterministic equivalence ([Cheng and Montanari, 2024](#); [Misiakiewicz and Saeed, 2024](#)): law of large samples/dimensions in random matrix theory

$$\text{Tr}(\mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda)^{-1}) \sim \text{Tr}(\Sigma (\Sigma + \lambda_*)^{-1}), w.h.p.$$

- \sim can be **asymptotic** or **non-asymptotic** at the rate of $\mathcal{O}(1/\sqrt{n})$.
- λ_* is the non-negative solution to the self-consistent equation
$$n - \frac{\lambda}{\lambda_*} = \text{Tr}(\Sigma (\Sigma + \lambda_*)^{-1}).$$

Our results

Theorem (Deterministic equivalence of estimator's norm)

We have a bias-variance decomposition $\mathbb{E}_\varepsilon \|\hat{\beta}\|_2^2 = \mathcal{B}_{\mathcal{N},\lambda} + \mathcal{V}_{\mathcal{N},\lambda}$.

For well-behaved data, we have

$$\mathcal{B}_{\mathcal{N},\lambda} := \langle \beta_*, \Sigma^2(\Sigma + \lambda_*)^{-2} \beta_* \rangle + \frac{\text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n} \frac{\lambda_*^2 \langle \beta_*, \Sigma(\Sigma + \lambda_*)^{-2} \beta_* \rangle}{1 - \frac{1}{n} \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})},$$

$$\mathcal{V}_{\mathcal{N},\lambda} := \frac{\sigma^2 \text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n - \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})}.$$

Remark: Which model capacity suffices to characterize the test risk?

- Norm-based capacity: ✓ ☺
- effective dimension-style $\text{Tr}(\Sigma(\Sigma + \lambda I)^{-1})$: ✗ ☺

Our results

Theorem (Deterministic equivalence of estimator's norm)

We have a bias-variance decomposition $\mathbb{E}_\varepsilon \|\hat{\beta}\|_2^2 = \mathcal{B}_{\mathcal{N},\lambda} + \mathcal{V}_{\mathcal{N},\lambda}$.

For well-behaved data, we have

$$\mathcal{B}_{\mathcal{N},\lambda} := \langle \beta_*, \Sigma^2(\Sigma + \lambda_*)^{-2} \beta_* \rangle + \frac{\text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n} \frac{\lambda_*^2 \langle \beta_*, \Sigma(\Sigma + \lambda_*)^{-2} \beta_* \rangle}{1 - \frac{1}{n} \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})} ,$$

$$\mathcal{V}_{\mathcal{N},\lambda} := \frac{\sigma^2 \text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n - \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})} .$$

Remark: Which model capacity suffices to characterize the test risk?

- Norm-based capacity: ✓ ☺
- effective dimension-style $\text{Tr}(\Sigma(\Sigma + \lambda I)^{-1})$: ✗ ☺

Our results

Theorem (Deterministic equivalence of estimator's norm)

We have a bias-variance decomposition $\mathbb{E}_\varepsilon \|\hat{\beta}\|_2^2 = \mathcal{B}_{\mathcal{N},\lambda} + \mathcal{V}_{\mathcal{N},\lambda}$.

For well-behaved data, we have

$$\mathcal{B}_{\mathcal{N},\lambda} := \langle \beta_*, \Sigma^2(\Sigma + \lambda_*)^{-2} \beta_* \rangle + \frac{\text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n} \frac{\lambda_*^2 \langle \beta_*, \Sigma(\Sigma + \lambda_*)^{-2} \beta_* \rangle}{1 - \frac{1}{n} \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})} ,$$

$$\mathcal{V}_{\mathcal{N},\lambda} := \frac{\sigma^2 \text{Tr}(\Sigma(\Sigma + \lambda_*)^{-2})}{n - \text{Tr}(\Sigma^2(\Sigma + \lambda_*)^{-2})} .$$

Remark: Which model capacity suffices to characterize the test risk?

- Norm-based capacity: ✓ 😊
- effective dimension-style $\text{Tr}(\Sigma(\Sigma + \lambda I)^{-1})$: ✗ 😞

Example: Relationship under isotropic features ($\Sigma = I_d$)

- Test risk R_λ and norm N_λ formulates a cubic curve (complex but precise).
 - min-norm interpolator ($\lambda = 0$):

$$R_0 = \begin{cases} N_0 - \|\beta_*\|_2^2; & \text{in under-parameterized regimes} \\ \sqrt{[N_0 - (\|\beta_*\|_2^2 - \sigma^2)]^2 + 4\|\beta_*\|_2^2\sigma^2} - \sigma^2. \end{cases}$$

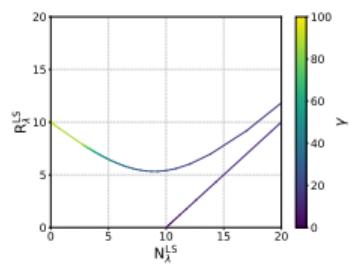
- optimal regularization $\lambda = \frac{d\sigma^2}{\|\beta^*\|_2^2}$ (Wu and Xu, 2020): $R_\lambda = \|\beta^*\|_2^2 - N_\lambda$

- $\lambda \rightarrow \infty$: $R_\lambda = (\|\beta^*\|_2 - \sqrt{N_\lambda})^2$

Example: Relationship under isotropic features ($\Sigma = I_d$)

- Test risk R_λ and norm N_λ formulates a cubic curve (complex but precise).
- min-norm interpolator ($\lambda = 0$):

$$R_0 = \begin{cases} N_0 - \|\beta_*\|_2^2; & \text{in under-parameterized regimes} \\ \sqrt{[N_0 - (\|\beta_*\|_2^2 - \sigma^2)]^2 + 4\|\beta_*\|_2^2\sigma^2} - \sigma^2. \end{cases}$$



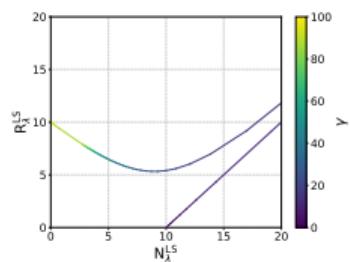
- optimal regularization $\lambda = \frac{d\sigma^2}{\|\beta_*\|_2^2}$ (Wu and Xu, 2020): $R_\lambda = \|\beta_*\|_2^2 - N_\lambda$

- $\lambda \rightarrow \infty$: $R_\lambda = (\|\beta_*\|_2 - \sqrt{N_\lambda})^2$

Example: Relationship under isotropic features ($\Sigma = I_d$)

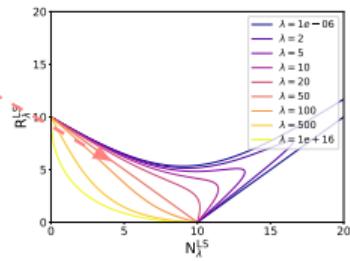
- Test risk R_λ and norm N_λ formulates a cubic curve (complex but precise).
- min-norm interpolator ($\lambda = 0$):

$$R_0 = \begin{cases} N_0 - \|\beta_*\|_2^2; & \text{in under-parameterized regimes} \\ \sqrt{[N_0 - (\|\beta_*\|_2^2 - \sigma^2)]^2 + 4\|\beta_*\|_2^2\sigma^2} - \sigma^2. \end{cases}$$



- optimal regularization $\lambda = \frac{d\sigma^2}{\|\beta^*\|_2^2}$ (Wu and Xu, 2020): $R_\lambda = \|\beta_*\|_2^2 - N_\lambda$

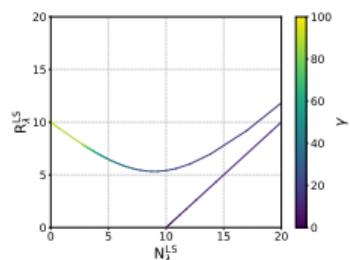
$\bullet \lambda \rightarrow \infty: R_\lambda = (\|\beta_*\|_2 - \sqrt{N_\lambda})^2$



Example: Relationship under isotropic features ($\Sigma = I_d$)

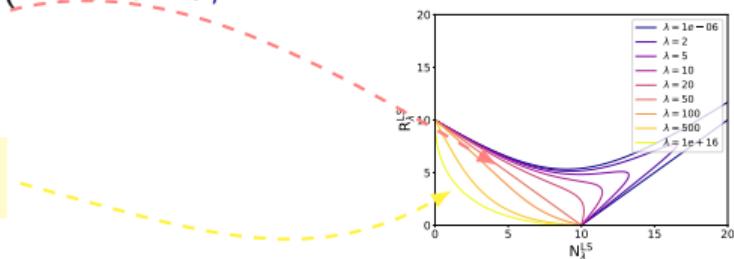
- Test risk R_λ and norm N_λ formulates a cubic curve (complex but precise).
- min-norm interpolator ($\lambda = 0$):

$$R_0 = \begin{cases} N_0 - \|\beta_*\|_2^2; & \text{in under-parameterized regimes} \\ \sqrt{[N_0 - (\|\beta_*\|_2^2 - \sigma^2)]^2 + 4\|\beta_*\|_2^2\sigma^2} - \sigma^2. \end{cases}$$



- optimal regularization $\lambda = \frac{d\sigma^2}{\|\beta^*\|_2^2}$ (Wu and Xu, 2020): $R_\lambda = \|\beta_*\|_2^2 - N_\lambda$

$$\bullet \lambda \rightarrow \infty: R_\lambda = (\|\beta_*\|_2 - \sqrt{N_\lambda})^2$$

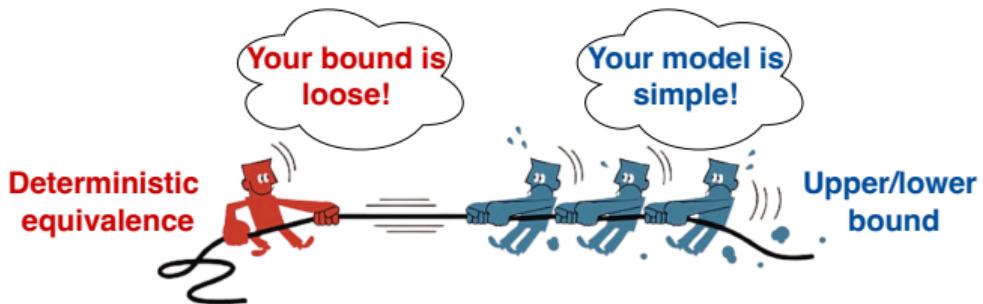


Precise analysis via deterministic equivalence

- Precisely describe the learning curve.
 - phase transitions, (non-)monotonicity, etc.
- Enables *accurate comparison* between estimators/algorithms.
 - **Foundations of scaling law**: data or parameter under the same budget, etc.

Precise analysis via deterministic equivalence

- Precisely describe the learning curve.
 - phase transitions, (non-)monotonicity, etc.
- Enables accurate *comparison* between estimators/algorithms.
 - Foundations of scaling law: data or parameter under the same budget, etc.



Which model capacity is suitable (for neural networks)?

Table 1: Complexity measures compared in the empirical study (Jiang et al., 2020), and their correlation with generalization.

name	definition	rank correlation
Parameter Frobenius norm	$\sum_{i=1}^L \ W_i\ _F^2$	0.073
Frobenius distance to initialization	$\sum_{i=1}^L \ W_i - W_i^0\ _F^2$	-0.263
Spectral complexity	$\prod_{i=1}^L \ W_i\ \left(\sum_{i=1}^L \frac{\ W_i\ _{2,1}^{3/2}}{\ W_i\ ^{3/2}} \right)^{2/3}$	-0.537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle W, \nabla_W \ell(h_W(x_i), y_i) \rangle$	0.078
Path-norm	$\sum_{(i_0, \dots, i_L)} \prod_{j=1}^L (W_{i_j, i_{j-1}})^2$	0.373

Which model capacity is suitable (for neural networks)?

Table 1: Complexity measures compared in the empirical study (Jiang et al., 2020), and their correlation with generalization.

name	definition	rank correlation
Parameter Frobenius norm	$\sum_{i=1}^L \ \mathbf{W}_i\ _F^2$	0.073
Frobenius distance to initialization	$\sum_{i=1}^L \ \mathbf{W}_i - \mathbf{W}_i^0\ _F^2$	-0.263
Spectral complexity	$\prod_{i=1}^L \ \mathbf{W}_i\ \left(\sum_{i=1}^L \frac{\ \mathbf{W}_i\ _{2,1}^{3/2}}{\ \mathbf{W}_i\ ^{3/2}} \right)^{2/3}$	-0.537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle \mathbf{W}, \nabla_{\mathbf{W}} \ell(h_{\mathbf{W}}(\mathbf{x}_i), y_i) \rangle$	0.078
Path-norm	$\sum_{(i_0, \dots, i_L)} \prod_{j=1}^L \left(\mathbf{W}_{i_j, i_{j-1}} \right)^2$	0.373

Which model capacity is suitable (for neural networks)?

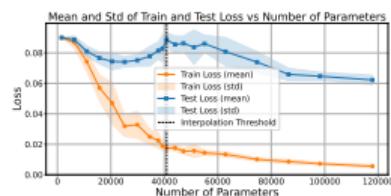
Table 1: Complexity measures compared in the empirical study (Jiang et al., 2020), and their correlation with generalization.

name	definition	rank correlation
Parameter Frobenius norm	$\sum_{i=1}^L \ \mathbf{W}_i\ _F^2$	0.073
Frobenius distance to initialization	$\sum_{i=1}^L \ \mathbf{W}_i - \mathbf{W}_i^0\ _F^2$	-0.263
Spectral complexity	$\prod_{i=1}^L \ \mathbf{W}_i\ \left(\sum_{i=1}^L \frac{\ \mathbf{W}_i\ _{2,1}^{3/2}}{\ \mathbf{W}_i\ ^{3/2}} \right)^{2/3}$	-0.537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle \mathbf{W}, \nabla_{\mathbf{W}} \ell(h_{\mathbf{W}}(\mathbf{x}_i), y_i) \rangle$	0.078
Path-norm	$\sum_{(i_0, \dots, i_L)} \prod_{j=1}^L (\mathbf{W}_{i_j, i_{j-1}})^2$	0.373

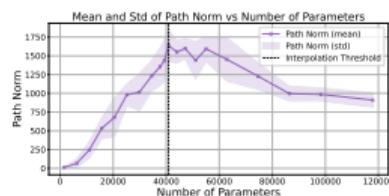
Which model capacity is suitable (for neural networks)?

Table 1: Complexity measures compared in the empirical study (Jiang et al., 2020), and their correlation with generalization.

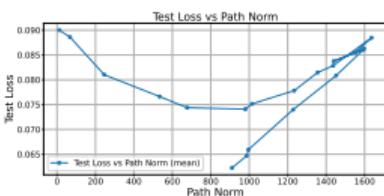
name	definition	rank correlation
Parameter Frobenius norm	$\sum_{i=1}^L \ \mathbf{W}_i\ _F^2$	0.073
Frobenius distance to initialization	$\sum_{i=1}^L \ \mathbf{W}_i - \mathbf{W}_i^0\ _F^2$	-0.263
Spectral complexity	$\prod_{i=1}^L \ \mathbf{W}_i\ \left(\sum_{i=1}^L \frac{\ \mathbf{W}_i\ _{2,1}^{3/2}}{\ \mathbf{W}_i\ ^{3/2}} \right)^{2/3}$	-0.537
Fisher-Rao	$\frac{(L+1)^2}{n} \sum_{i=1}^n \langle \mathbf{W}, \nabla_{\mathbf{W}} \ell(h_{\mathbf{W}}(\mathbf{x}_i), y_i) \rangle$	0.078
Path-norm	$\sum_{(i_0, \dots, i_L)} \prod_{j=1}^L (\mathbf{W}_{i_j, i_{j-1}})^2$	0.373



(a) Test (training) Loss vs. p



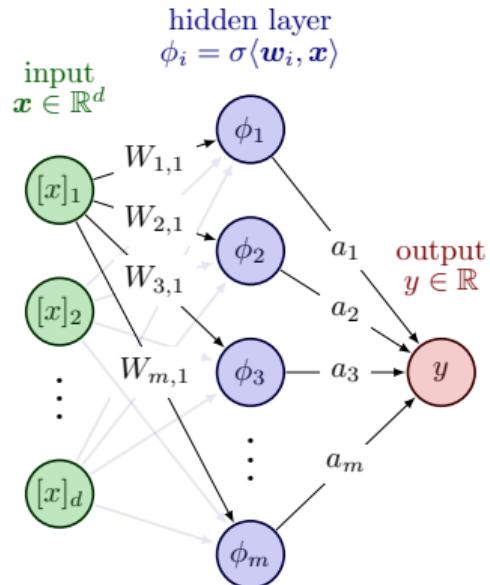
(b) Path-norm vs. p



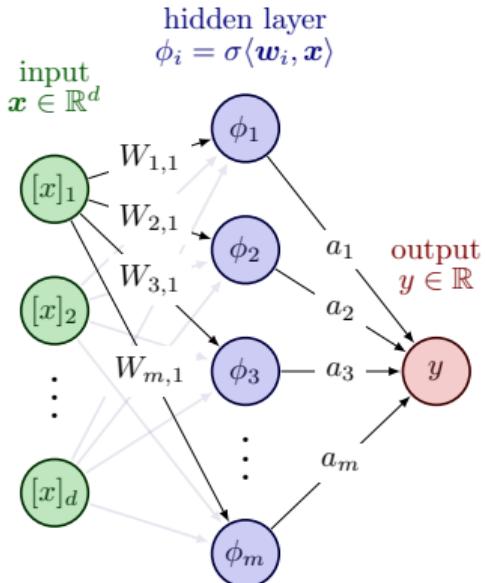
(c) Test Loss vs. Path-norm

Figure 5: Experiments on two-layer neural networks.

Two-layer neural networks, path norm



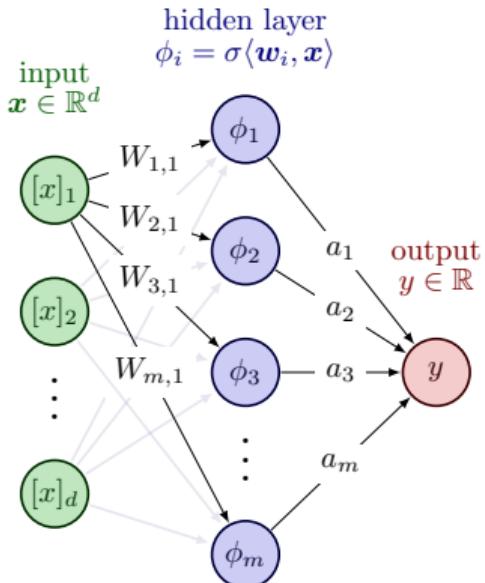
Two-layer neural networks, path norm



ℓ_1 -path norm (Neyshabur et al., 2015)

$$\|\boldsymbol{\theta}\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$$

Two-layer neural networks, path norm



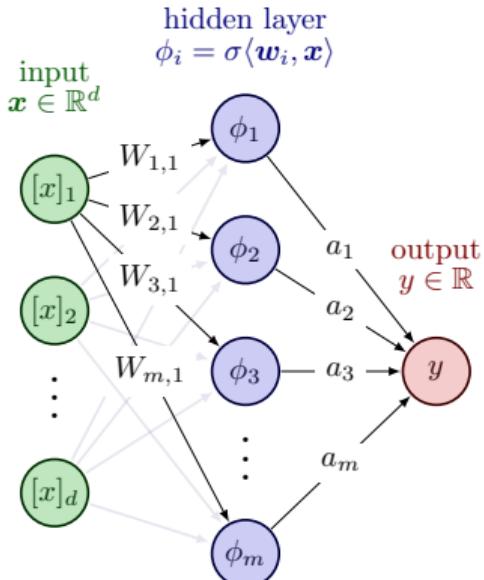
ℓ_1 -path norm (Neyshabur et al., 2015)

$$\|\boldsymbol{\theta}\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$$

- equivalent to Barron spaces \mathcal{B}
(Barron, 1993; E et al., 2021)

$$\mathcal{B} := \cup_{\mu \in \mathcal{P}(\mathcal{W})} \{f_a : \|a\|_{L^2(\mu)} < \infty\}$$

Two-layer neural networks, path norm



ℓ_1 -path norm (Neyshabur et al., 2015)

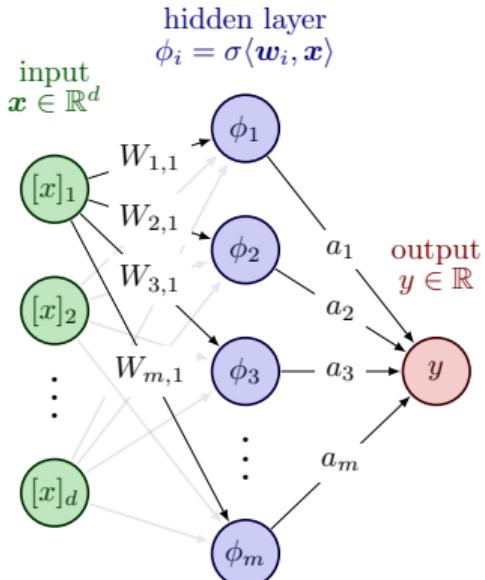
$$\|\boldsymbol{\theta}\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$$

- equivalent to Barron spaces \mathcal{B} (Barron, 1993; E et al., 2021)

$$\mathcal{B} := \cup_{\mu \in \mathcal{P}(\mathcal{W})} \{f_a : \|a\|_{L^2(\mu)} < \infty\}$$

- Variation in only a few directions (Parhi and Nowak, 2022)

Two-layer neural networks, path norm



ℓ_1 -path norm (Neyshabur et al., 2015)

$$\|\boldsymbol{\theta}\|_{\mathcal{P}} := \frac{1}{m} \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1$$

- equivalent to Barron spaces \mathcal{B} (Barron, 1993; E et al., 2021)

$$\mathcal{B} := \cup_{\mu \in \mathcal{P}(\mathcal{W})} \{f_a : \|a\|_{L^2(\mu)} < \infty\}$$

- Variation in only a few directions (Parhi and Nowak, 2022)

Can neural networks identify this structure?



Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}



Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}



Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}

No **Curse of Dimensionality**: NNs adapt to directional smoothness.



Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}

No **Curse of Dimensionality**: NNs adapt to directional smoothness.

□ Track sample complexity (via metric entropy) and dimension dependence





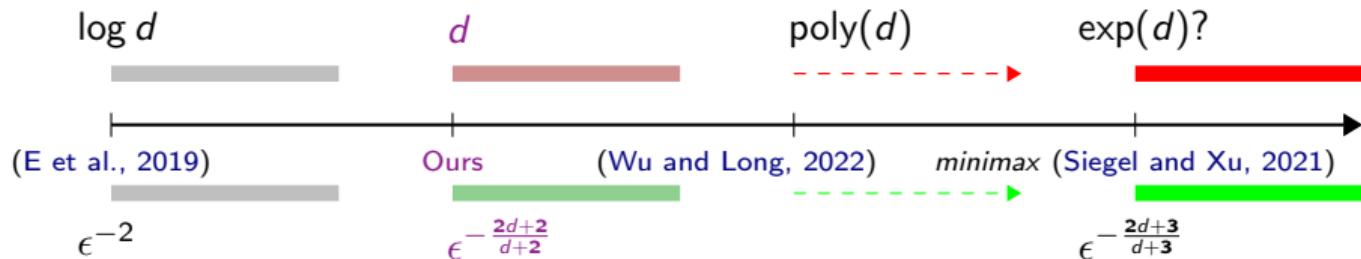
Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}

No **Curse of Dimensionality**: NNs adapt to directional smoothness.

□ Track sample complexity (via metric entropy) and dimension dependence





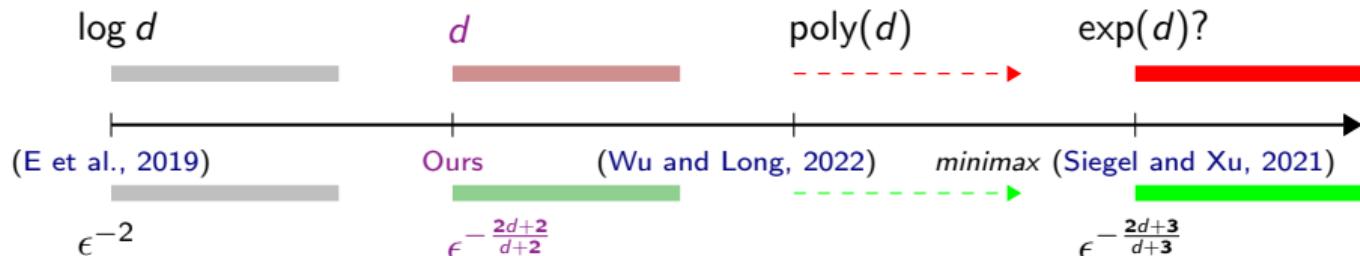
Theorem (Informal, sample complexity of learning $f^* \in \mathcal{B}$)

To achieve ϵ -excess risk,

- Kernel methods require $\Omega(\epsilon^{-d})$ samples.
- Two-layer neural networks require $\Omega(\epsilon^{-\frac{2d+2}{d+2}})$ samples. smaller than ϵ^{-2}

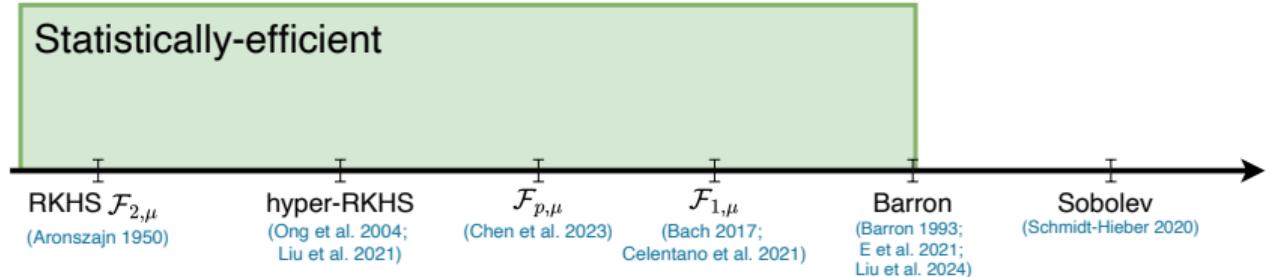
No **Curse of Dimensionality**: NNs adapt to directional smoothness.

□ Track sample complexity (via metric entropy) and dimension dependence

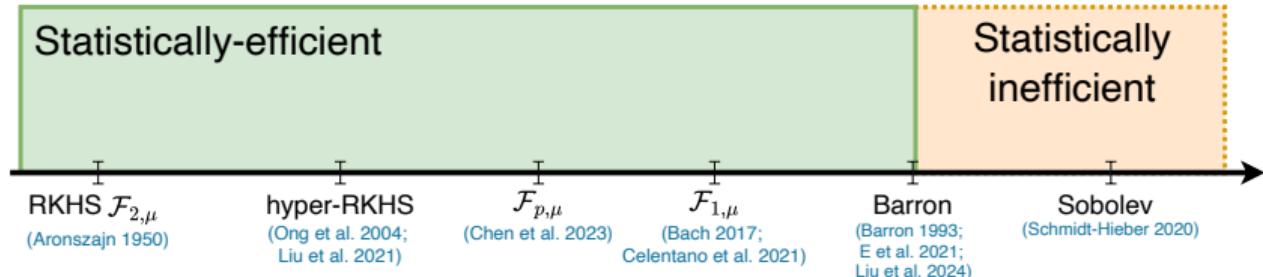


The “best” trade-off between ϵ and d .

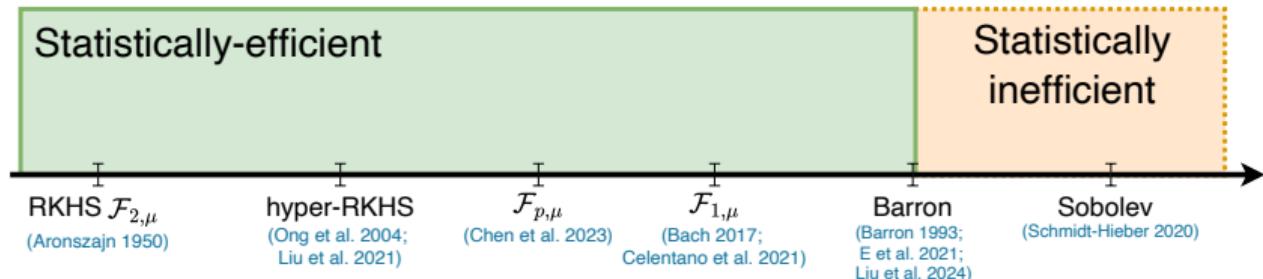
Which function class can be efficiently learned by neural networks



Which function class can be efficiently learned by neural networks

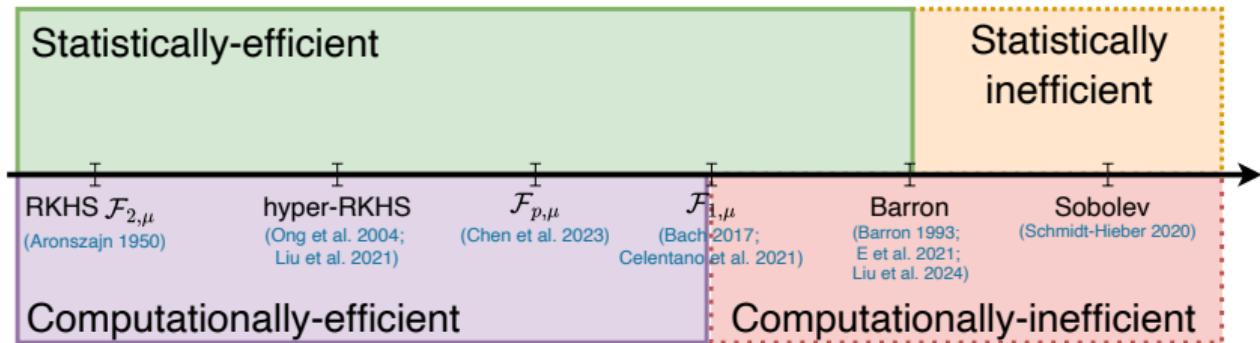


Which function class can be efficiently learned by neural networks

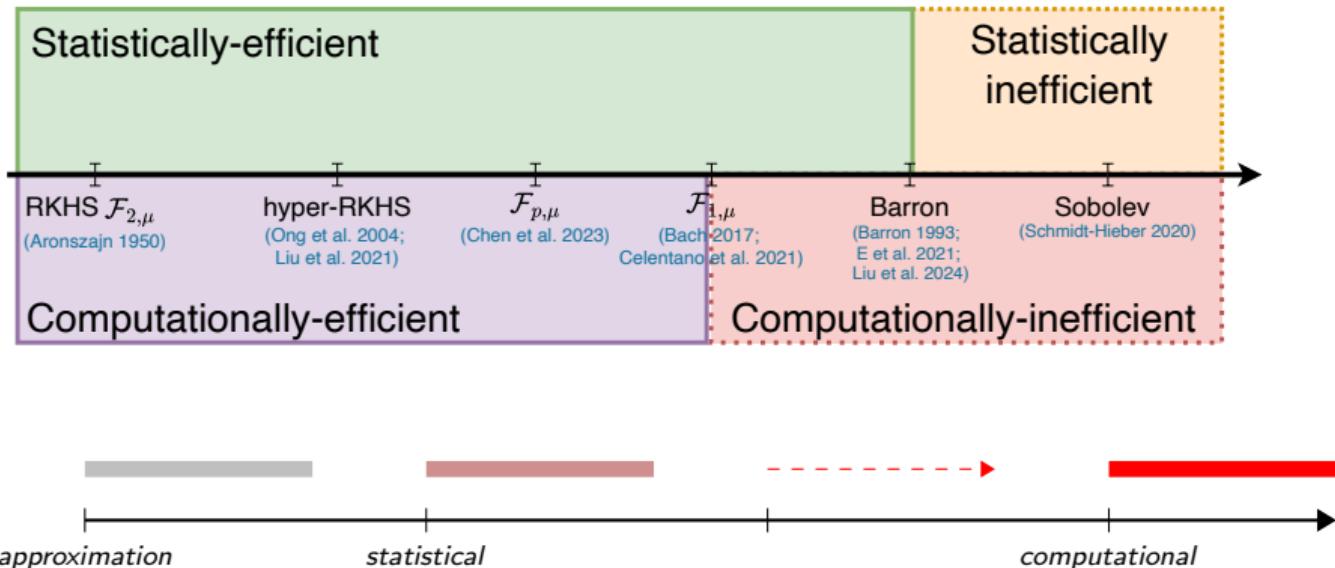


Optimization in Barron spaces is NP hard: curse of dimensionality!
(Bach, 2017)

Which function class can be efficiently learned by neural networks



Which function class can be efficiently learned by neural networks

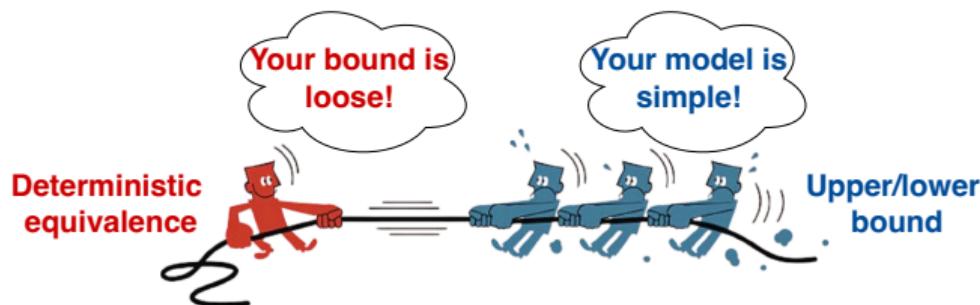


- ReLU neurons (Chen and Narayanan, 2023)
- Low-dimensional polynomials (Arous et al., 2021; Lee et al., 2024)

Deep learning phenomena \Rightarrow interesting mathematical problems

Be aware of model capacity!

- Reshape bias-variance trade-offs, double descent, scaling law under proper ℓ_2 norm-based capacity via **deterministic equivalence**.

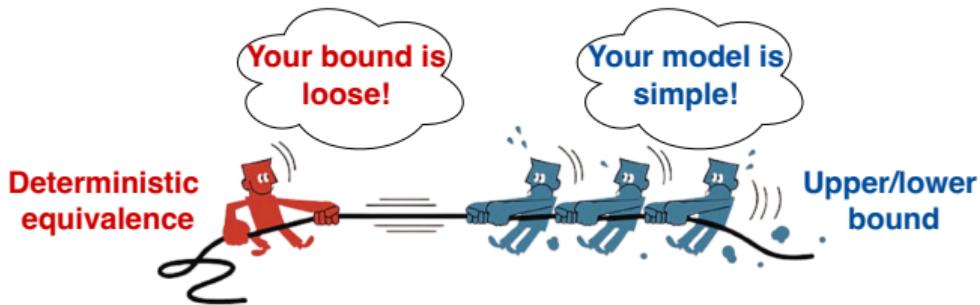


Takeaway messages

Deep learning phenomena \Rightarrow interesting mathematical problems

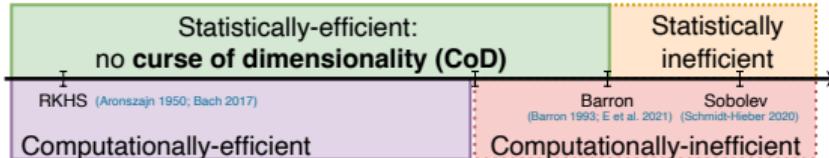
Be aware of model capacity!

- Reshape bias-variance trade-offs, double descent, scaling law under proper ℓ_2 norm-based capacity via **deterministic equivalence**.



Which function class can be **efficiently** learned by neural networks?

- Neural networks can adapt to low-dimensional structure and avoid CoD!

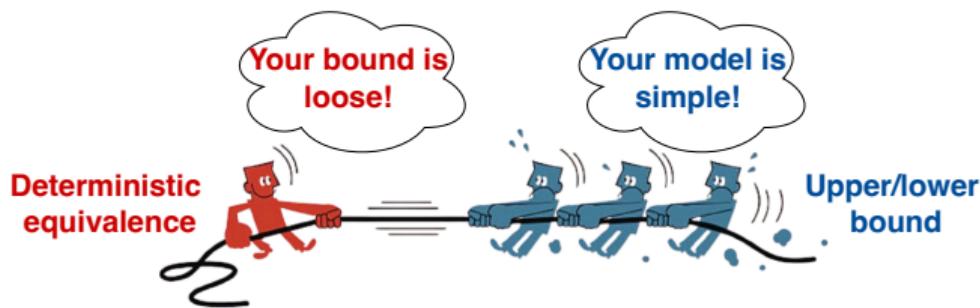


Takeaway messages

Deep learning phenomena \Rightarrow interesting mathematical problems

Be aware of model capacity!

- Reshape bias-variance trade-offs, double descent, scaling law under proper ℓ_2 norm-based capacity via **deterministic equivalence**.



Which function class can be **efficiently** learned by neural networks?

- Neural networks can adapt to low-dimensional structure and avoid CoD!

Theoretical advances \Rightarrow principled guidance in practical problems

How does our theory contribute to practical fine-tuning problems?

- One-step full gradient can be sufficient! [\[GitHub\]](#)

References

- Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- Francis Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(1):629–681, 2017.
- Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory*, 39(3):930–945, 1993.

- Peter Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, 1998.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *the National Academy of Sciences*, 116(32):15849–15854, 2019.
- Sitan Chen and Shyam Narayanan. A faster and simpler algorithm for learning shallow networks. *arXiv preprint arXiv:2307.12496*, 2023.
- Chen Cheng and Andrea Montanari. Dimension free ridge regression. *The Annals of Statistics*, 52(6):2879–2912, 2024.
- Weinan E, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-layer neural networks. *Communications in Mathematical Sciences*, 17(5):1407–1425, 2019.

- Weinan E, Chao Ma, and Lei Wu. The barron space and the flow-induced function spaces for neural network models. *Constructive Approximation*, pages 1–38, 2021.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. *Neural computation*, 4(1):1–58, 1992.
- Per Christian Hansen. Analysis of discrete ill-posed problems by means of the l-curve. *SIAM Review*, 34(4):561–580, 1992.
- Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.

- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani.
Surprises in high-dimensional ridgeless least squares interpolation. *Annals of Statistics*, 50(2):949–986, 2022.
- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Jason D Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional polynomials with sgd near the information-theoretic limit. *arXiv preprint arXiv:2406.01581*, 2024.

- Tengyuan Liang, Alexander Rakhlin, and Xiyu Zhai. On the multiple descent of minimum-norm interpolants and restricted lower isometry of kernels. In *Conference on Learning Theory*, pages 2683–2711, 2020.
- Theodor Misiakiewicz and Basil Saeed. A non-asymptotic theory of kernel ridge regression: deterministic equivalents, test error, and gcv estimator. *arXiv preprint arXiv:2403.08938*, 2024.
- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. In *International Conference on Learning Representations*, 2017.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*, 2019.

- Brady Neal. On the bias-variance tradeoff: Textbooks need an update. *arXiv preprint arXiv:1912.08286*, 2019.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In *Conference on Learning Theory*, pages 1376–1401. PMLR, 2015.
- Andrew Ng and Tengyu Ma. CS229 lecture notes. 2023. URL https://cs229.stanford.edu/main_notes.pdf.
- Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+3 phases of compute-optimal neural scaling laws. *arXiv preprint arXiv:2405.15074*, 2024.
- Rahul Parhi and Robert D Nowak. Near-minimax optimal estimation with shallow ReLU neural networks. *IEEE Transactions on Information Theory*, 2022.

- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pages 1177–1184, 2007.
- Pedro Savarese, Itay Evron, Daniel Soudry, and Nathan Srebro. How do infinite width bounded norm networks look in function space? In *Conference on Learning Theory*, pages 2667–2690. PMLR, 2019.
- Jonathan W Siegel and Jinchao Xu. Sharp bounds on the approximation rates, metric entropy, and n -widths of shallow neural networks. *arXiv preprint arXiv:2101.12365*, 2021.
- Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- Andrew Gordon Wilson. Deep learning is not so mysterious or different. *arXiv preprint arXiv:2503.02113*, 2025.

- Denny Wu and Ji Xu. On the optimal weighted ℓ_2 regularization in overparameterized linear regression. In *Advances in Neural Information Processing Systems*, pages 10112–10123, 2020.
- Lei Wu and Jihao Long. A spectral-based analysis of the separation between two-layer neural networks and linear methods. *Journal of Machine Learning Research*, 119:1–34, 2022.